

MENU

SEARCH

INDEX

DETAIL

E5388

1/1



JAPANESE PATENT OFFICE

PATENT ABSTRACTS OF JAPAN

(11)Publication number: 2000115243

(43)Date of publication of application: 21.04.2000

(51)Int.Cl.

H04L 12/56
G06F 13/00
H04L 12/46
H04L 12/28
H04L 12/18
H04L 13/08

(21)Application number: 10282205

(71)Applicant:

HITACHI LTD
HITACHI INFORMATION
TECHNOLOGY CO LTD

(22)Date of filing: 05.10.1998

(72)Inventor:

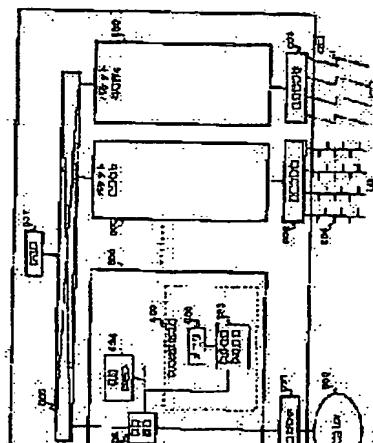
AKAHA SHINICHI
SUKAI KAZUO
AIMOTO TAKESHI
MATSUYAMA NOBUHITO
SAKO YOSHITO
SEKINO HIROSHI

(54) PACKET REPEATER AND MULTI-CAST SPEEDUP SYSTEM

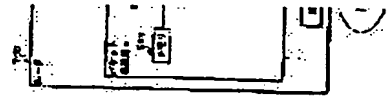
(57)Abstract:

PROBLEM TO BE SOLVED: To obtain a system that retrieves a multi-cast route for the retrieval of transmission port information used to transfer a received multi-cast packet succeedingly at a high speed.

SOLUTION: A (p -th power of 2) tree is configured, by forming a route address and connecting a destination address and a transmission source address in this order, using nodes in total of (p -th power of 2 minus one), being the sum of one node of a binary retrieval tree corresponding to a bit pattern of the route address and nodes for p -stages directly beneath the node, and is stored in a memory 300, not one bit but p -



consecutive bits of the route address obtained by connecting a destination address and a transmission source address to each other in this order of a received multi-cast packet are checked at the same time and the retrieval tree stored in the memory 300 is retrieved by using consecutive bit values for key so as to complete retrieval processing and passing through nodes by number of $(\text{bit number of retrieval keys} \div p)$ times independently of number of entries.



LEGAL STATUS

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

Copyright (C); 1998 Japanese Patent Office

[MENU](#)

[SEARCH](#)

[INDEX](#)

[DETAIL](#)

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開2000-115243

(P2000-115243A)

(43) 公開日 平成12年4月21日 (2000.4.21)

(51) Int.Cl. ⁷	識別記号	F I	テマート* (参考)
H 0 4 L 12/56		H 0 4 L 11/20	1 0 2 Z 5 B 0 8 9
G 0 6 F 13/00	3 5 5	G 0 6 F 13/00	3 5 5 5 K 0 3 0
H 0 4 L 12/46		H 0 4 L 13/08	5 K 0 3 3
12/28		11/00	3 1 0 C 5 K 0 3 4
12/18		11/18	

審査請求 未請求 請求項の枚数 23 O L (全 37 頁) 最終頁に続く

(21) 出願番号 特願平10-282205

(22) 出願日 平成10年10月5日 (1998.10.5)

(71) 出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(71) 出願人 000153454

株式会社日立インフォメーションテクノロ
ジー

神奈川県足柄上郡中井町境456番地

(72) 発明者 赤羽 真一

東京都国分寺市東恋ヶ窪一丁目280番地

株式会社日立製作所中央研究所内

(74) 代理人 100068504

弁理士 小川 勝男

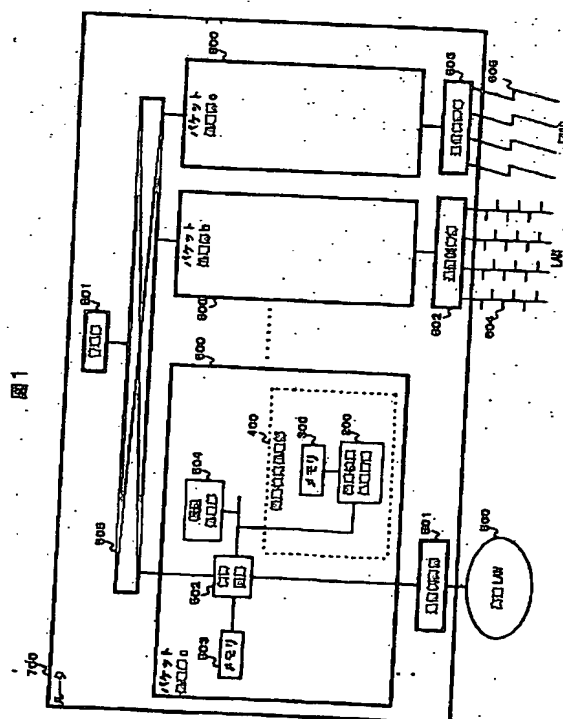
最終頁に続く

(54) 【発明の名称】 パケット中継装置およびマルチキャスト高速化方式

(57) 【要約】

【課題】 受信したマルチキャストパケットを次に転送すべき送信ポート情報を検索するマルチキャスト経路検索を高速に行う方式を提案する。

【解決手段】 宛先アドレスと送信元アドレスをこの順に連結して経路アドレスとし、上記経路アドレスのビットパターンに対応した2分木検索木の一つのノードとその直下のp段分のノードの計(2のp乗-1)個のノードを一つの2のp乗分木ノードにして、2のp乗分木を構成してメモリ300に格納し、受信したマルチキャストパケットの宛先アドレスと送信元アドレスをこの順に連結した経路アドレスを、一つのノードで1ビットではなく、連続するpビットを同時に検査し、上記連続するビットの値によってメモリ300に格納された検索木を検索することにより、エントリ数に依存せず、かつ、最大、(検索のキーのビット数÷p)回ノードを渡るだけで検索処理を終了することができる。



【特許請求の範囲】

【請求項1】受信したマルチキャストパケットの宛先アドレスと送信元アドレスを検索のキーにして上記マルチキャストパケットに対する経路情報を検索し、上記経路情報に従い1つ又は複数のポートに上記マルチキャストパケットを転送するパケット中継装置において、宛先アドレスと送信元アドレスとを連結した経路アドレスのビットパターンに対応した2分木構造の検索木情報を格納する記憶手段と、

受信したマルチキャストパケットから宛先アドレスと送信元アドレスとを読み出し、その読み出した宛先アドレスと送信元アドレスとを連結して経路アドレスを生成し、その生成された経路アドレスと上記記憶手段に格納されている経路アドレスとの一致判定を行う一致判定手段とを有することを特徴とするパケット中継装置。

【請求項2】受信したマルチキャストパケットの宛先アドレスと送信元アドレスを検索のキーにして上記マルチキャストパケットに対する経路情報を検索し、上記経路情報に従い1つ又は複数のポートに上記マルチキャストパケットを転送するパケット中継装置において、宛先アドレスと送信元アドレスとをこの順に連結した経路アドレスのビットパターンに対応した2分木構造の検索木情報を格納する記憶手段と、

受信したマルチキャストパケットから宛先アドレスと送信元アドレスとを読み出し、その読み出した宛先アドレスと送信元アドレスとをこの順に連結して経路アドレスを生成し、その生成された経路アドレスと上記記憶手段に格納されている経路アドレスとの一致判定を行う一致判定手段とを有することを特徴とするパケット中継装置。

【請求項3】上記一致判定手段は、上記受信したマルチキャストパケットから生成した経路アドレスの上位ビットから1ビットづつ、上記記憶手段に格納されている経路アドレスとの一致判定を行うことを特徴とする請求項1又は請求項2の何れかに記載のパケット中継装置。

【請求項4】受信したマルチキャストパケットの宛先アドレスと送信元アドレスを検索のキーにして上記マルチキャストパケットに対する経路情報を検索し、上記経路情報に従い1つ又は複数のポートに上記マルチキャストパケットを転送するパケット中継装置において、宛先アドレスと送信元アドレスとを連結した経路アドレスのビットパターンに対応した2分木構造の検索木情報を格納する記憶手段と、

受信したマルチキャストパケットから宛先アドレスと送信元アドレスとを読み出し、その読み出した宛先アドレスと上記記憶手段に格納されている経路アドレスとの一致判定を行い、その読み出した送信元アドレスと上記記憶手段に格納されている経路アドレスとの一致判定を行う一致判定手段とを有することを特徴とするパケット中継装置。

【請求項5】上記一致判定手段は、受信したマルチキャストパケットから読み出した宛先アドレスと上記記憶手段に格納されている経路アドレスとの一致判定を、受信したマルチキャストパケットから読み出した送信元アドレスと上記記憶手段に格納されている経路アドレスとの一致判定より先に行うことを特徴とする請求項4に記載のパケット中継装置。

【請求項6】受信したマルチキャストパケットの宛先アドレスと送信元アドレスを検索のキーにして上記マルチキャストパケットに対する経路情報を検索し、上記経路情報に従い1つ又は複数のポートに上記マルチキャストパケットを転送するパケット中継装置において、宛先アドレスと送信元アドレスとを連結した経路アドレスのビットパターンに対応した2の p 乗分木構造（ p は自然数）の検索木情報を格納する記憶手段と、

受信したマルチキャストパケットから宛先アドレスと送信元アドレスとを読み出し、その読み出した宛先アドレスと送信元アドレスとを連結して経路アドレスを生成し、その生成された経路アドレスと上記記憶手段に格納されている経路アドレスとの一致判定を行う一致判定手段とを有することを特徴とするパケット中継装置。

【請求項7】受信したマルチキャストパケットの宛先アドレスと送信元アドレスを検索のキーにして上記マルチキャストパケットに対する経路情報を検索し、上記経路情報に従い1つ又は複数のポートに上記マルチキャストパケットを転送するパケット中継装置において、宛先アドレスと送信元アドレスとをこの順に連結した経路アドレスのビットパターンに対応した2の p 乗分木構造（ p は自然数）の検索木情報を格納する記憶手段と、受信したマルチキャストパケットから宛先アドレスと送信元アドレスとを読み出し、その読み出した宛先アドレスと送信元アドレスとをこの順に連結して経路アドレスを生成し、その生成された経路アドレスと上記記憶手段に格納されている経路アドレスとの一致判定を行う一致判定手段とを有することを特徴とするパケット中継装置。

【請求項8】上記一致判定手段は、上記受信したマルチキャストパケットから生成した経路アドレスの上位ビットから p ビットづつ、上記記憶手段に格納されている経路アドレスとの一致判定を行うことを特徴とする請求項6又は請求項7の何れかに記載のパケット中継装置。

【請求項9】受信したマルチキャストパケットの宛先アドレスと送信元アドレスを検索のキーにして上記マルチキャストパケットに対する経路情報を検索し、上記経路情報に従い1つ又は複数のポートに上記マルチキャストパケットを転送するパケット中継装置において、宛先アドレスと送信元アドレスとを連結した経路アドレスのビットパターンに対応した2の p 乗分木構造（ p は自然数）の検索木情報を格納する記憶手段と、受信したマルチキャストパケットから宛先アドレスと送

信元アドレスとを読み出し、その読み出した宛先アドレスと上記記憶手段に格納されている経路アドレスとの一致判定を行い、その読み出した送信元アドレスと上記記憶手段に格納されている経路アドレスとの一致判定を行う一致判定手段とを有することを特徴とするバケット中継装置。

【請求項10】上記一致判定手段は、受信したマルチキャストバケットから読み出した宛先アドレスと上記記憶手段に格納されている経路アドレスとの一致判定を、受信したマルチキャストバケットから読み出した送信元アドレスと上記記憶手段に格納されている経路アドレスとの一致判定より先に行うことを特徴とする請求項9に記載のバケット中継装置。

【請求項11】複数のネットワークを相互に接続し、バケットを中継するバケット中継装置であって、受信したマルチキャストバケットの宛先アドレスと送信元アドレスを検索のキーにして上記マルチキャストバケットに対する経路情報を検索し、上記経路情報に従い1つ又は複数のポートに上記マルチキャストバケットを転送するバケット中継装置において、宛先アドレスと送信元アドレスを連結して経路アドレスとし、上記経路アドレスのビットパターンに対応した2分木構造に構成された検索木情報を格納する記憶手段と、受信したマルチキャストバケットの宛先アドレスと送信元アドレスを連結した経路アドレスの上位ビットから1ビットずつ検査し、その検査したビットの値によって上記記憶手段に格納された2分木構造をもつ検索木を検索する回路とを有することを特徴とするバケット中継装置。

【請求項12】複数のネットワークを相互に接続し、バケットを中継するバケット中継装置であって、受信したマルチキャストバケットの宛先アドレスと送信元アドレスを検索のキーにして上記マルチキャストバケットに対する経路情報を検索し、上記経路情報に従い1つ又は複数のポートに上記マルチキャストバケットを転送するバケット中継装置において、送信元アドレスと宛先アドレスとを、宛先アドレス、送信元アドレスの順に連結して経路アドレスとし、上記経路アドレスのビットパターンに対応した2分木構造に構成された検索木情報を格納する記憶手段と、受信したマルチキャストバケットの宛先アドレスと送信元アドレスをこの順に連結した経路アドレスの上位ビットから1ビットずつ検査し、その検査したビットの値によって上記記憶手段に格納された2分木構造をもつ検索木を検索する回路とを有することを特徴とするバケット中継装置。

【請求項13】複数のネットワークを相互に接続し、バケットを中継するバケット中継装置であって、受信したマルチキャストバケットの宛先アドレスと送信元アドレ

スを検索のキーにして上記マルチキャストバケットに対する経路情報を検索し、上記経路情報に従い1つ又は複数のポートに上記マルチキャストバケットを転送するバケット中継装置において、

送信元アドレスと宛先アドレスとを、宛先アドレス、送信元アドレスの順に連結して経路アドレスとし、上記経路アドレスのビットパターンに対応した2の p 乗分木構造に構成された検索情報を格納する記憶手段と、受信したマルチキャストバケットの宛先アドレスと送信元アドレスを連結した経路アドレスの上位ビットから p ビットずつ検査し、その検査したビットの値によって上記記憶手段に記憶された2の p 乗分木構造をもつ検索木を検索する回路とを有することを特徴とするバケット中継装置。

【請求項14】複数のネットワークを相互に接続し、バケットを中継するバケット中継装置であって、受信したマルチキャストバケットの宛先アドレスと送信元アドレスを検索のキーとして上記マルチキャストバケットに対する経路情報を検索し、上記経路情報に従い1つ又は複数のポートに上記マルチキャストバケットを転送するバケット中継装置において、送信元アドレスと宛先アドレスとをこの順に連結した経路アドレスのビットパターンに対応する2分木構造をとる検索テーブルを格納する記憶手段と、経路アドレスの上位 m ビットに対応する2の m 乗個の2分木のノードを、記憶手段上の決まった位置に展開し、それぞれのノードを、経路アドレスの第0ビットから第 $(m-1)$ ビットまでが取りうる値に1対1に対応させ、検索時には、経路アドレスの第0ビットから第 $(m-1)$ ビットの値に従って、展開されたノードの一つを選択し、第 m ビット以降は経路アドレスを1ビットずつ検査し、その検査ビットの値によって上記記憶手段に記憶された2分木構造をもつ検索木の枝を検索する回路とを有することを特徴とするバケット中継装置。

【請求項15】複数のネットワークを相互に接続し、バケットを中継するバケット中継装置であって、受信したマルチキャストバケットの宛先アドレスと送信元アドレスを検索のキーとして上記マルチキャストバケットに対する経路情報を検索し、上記経路情報に従い1つ又は複数のポートに上記マルチキャストバケットを転送するバケット中継装置において、送信元アドレスと宛先アドレスを、宛先アドレス、送信元アドレスの順に連結して経路アドレスとし、この経路アドレスのビットパターンに対応する2の p 乗分木構造をとる検索テーブルを格納する記憶手段と、経路アドレスの上位 m ビットに対応する2の m 乗個の2の p 乗分木のノードを、記憶手段上の決まった位置に展開し、それぞれのノードを、経路アドレスの第0ビットから第 $(m-1)$ ビットまでが取りうる値に1対1に対応させ、検索時には、経路アドレスの第0ビットから第

($m-1$)ビットの値に従って、展開されたノードの一つを選択し、第 m ビット以降は経路アドレスを p ビットづつ検査し、その検査したビットの値によって上記記憶手段に記憶された2の p 乗分木構造をもつ検索木の枝を検索する回路を有することを特徴とするバケット中継装置。

【請求項16】請求項12に記載のバケット中継装置において、

検索木情報を、1つの2分木ノードと、その直下につながる $p-1$ 段分の合計(2の p 乗 -1)個分の2分木ノードを一つの2の p 乗分木ノードにまとめ、まとめられる最下段の2の($p-1$)乗個の2分木ノードに、それより上段のノードに割り付けられていたエントリデータを埋め込み、2の p 乗分木ノードを、2分木を2の($p-1$)乗個分併せた形で構成することを特徴とするバケット中継装置。

【請求項17】請求項12に記載のバケット中継装置において、一つの2分木ノードと、その直下につながる $p-1$ 段分の合計(2の p 乗 -1)個分の2分木ノードを一つの2の p 乗分木ノードにまとめ、まとめられる最下段の2の($p-1$)乗個の2分木ノードに、それより上段のノードに割り付けられていたエントリデータを埋め込むことにより、2の p 乗分木ノードを、2分木を2の($p-1$)乗個分併せた形で構成することに加え、2分木を複数個併せるときに、一つだけ持てば良い要素を一つだけ保持するようにすることを特徴とする特徴とするバケット中継装置。

【請求項18】請求項15又は請求項16の何れかに記載のバケット中継装置において、

2の p 乗分木ノードを検索のために読むときにノード全てを読まずに、2の p 乗分木ノードを作成するときに併せた2の($p-1$)乗個の2分木ノードの内の、いずれか一つに対応するデータのみを読むようにすることを特徴とするバケット中継装置。

【請求項19】請求項17に記載のバケット中継装置において、

各ノードにそのノード自身のマスク長を格納するのではなく、そのノードの直ぐ下に繋がるノードのマスク長を格納することにより、ノードのマスク長を、そのノードのデータを読む前に知り、経路アドレスの、ノードのマスク長で示されるビット位置から、そのビット位置 $+p-1$ までの値に従い、ノードのデータの内、読み込む部分を選択することを特徴とするバケット中継装置。

【請求項20】ルータに入ってきたバケットの宛先アドレスから、このバケットを次に出力すべきポート、及び、次に送るべきルータ、又はホストのアドレスを高速に検索するための検索手段に関し、

検索を、宛先アドレスの上位ビットから1ビットづつ検査してゆく2分木検索により行い、検査を行うビット位置を、マスク長に対応させることにより、マスク付きの

一致検索を行い、検索木の各ノードの最初に読み込むデータ内に、そのノードにエントリが割り付けられているか否かを示すフラグを設け、最初に、このフラグを読み込み、エントリが割り付けられていないノードでは、経路情報を読み込まないようにすることにより、データの読み込み時間を短縮し、検索処理の高速化を図ることを特徴とするバケット中継装置。

【請求項21】請求項11乃至請求項13の何れかに記載のバケット中継装置において、

マルチキャストバケットの宛先アドレスの上位数ビットがマルチキャストバケットであることを示す識別ビットとして使用される場合、宛先アドレスから上記識別ビットを除いたものと送信元アドレスを連結して経路アドレスとすることを特徴とするバケット中継装置。

【請求項22】請求項12に記載のバケット中継装置において、

検索結果である送信ポート情報を検索木データから分離して記憶手段の別領域に格納し、検索木データには上記の分離した送信ポート情報が格納されているアドレスを保持することを特徴とする特徴とするバケット中継装置。

【請求項23】マルチキャストバケットのヘッダ情報を入力され、上記マルチキャストバケットに対する送信ポート等の転送情報を検索して出力する専用LSIにおいて、マルチキャストバケットの宛先アドレスと送信元アドレスをこの順に連結した経路アドレスの上位ビットから p ビットづつ検査し、上記検査ビットの値によって記憶手段に記憶された2の p 乗分木構造をもつ検索木を検索する回路を有することを特徴とする専用LSI。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、複数のネットワークを相互に接続し、バケットを中継するバケット中継装置、特にマルチキャストバケットの次転送先検索方式に関する。

【0002】

【従来の技術】ユーザの増加に伴い、インターネットを流れるトラフィック(バケット)が急増し、インターネットが大規模化・高速化している。また、現在のインターネット(Internet Protocolによるバケット通信網:以下IPネットワーク)では、従来のデータ系の通信だけでなく、インターネット電話やインターネット放送などのリアルタイム系アプリケーション(用途)も登場し、音声通信機能の取り込みや、放送機能の取り込みが進んでいる。このような状況下において、IPネットワークにおけるIPマルチキャスト技術は、インターネットにおける動画や音声やコンテンツなどのマルチメディアデータの配信に有効な技術として期待されている。また、IPネットワークを構成するバケット中継装置(ルータ)におけるIPマルチキャスト技術のサポート、IPマルチキ

キャスト技術の高速化が課題となっている。

【0003】ルータは、ある端末から特定の一つの端末にパケットを送信するユニキャスト通信の場合、受信パケットのヘッダ内の宛先IPアドレスに対応する経路情報(次に転送すべきルータあるいは端末のIPアドレス、およびルータ内の送信ポート番号)をルーティングテーブルから検索し、パケットを送信する。以下、ユニキャストの経路検索について簡単に説明する。

【0004】ルーティングテーブルは、宛先IPアドレスに対応するサブネットワークアドレスとサブネットマスク長および、上記の経路情報から構成される情報群を複数保持している。以下では、これらの情報群をエントリと呼ぶ。ここでサブネットワークとは、例えば宛先の企業網の様な端末の部分集合を示し、サブネットワークアドレスは部分集合のIPアドレスをアドレス情報として結合したアドレス群である。また、サブネットマスク長はIPアドレスのうち上位何ビットが上記のサブネットワークの識別子であるかを示す値である。ルータは、受信パケットの宛先IPアドレスにエントリ内のサブネットマスク長だけ上位ビットから有効なマスクをかけたものとエントリのサブネットワークアドレスを比較して、一致するエントリの経路情報を検索結果とする。このようにサブネットワーク単位でエントリを構成することによりルーティングテーブルのエントリ数を大幅に縮小し、検索処理の効率化を図っている。複数のエントリに一致した場合は、サブネットマスク長が最長のエントリの経路情報を検索結果とする。以下では、この検索方式を最長一致検索と呼ぶ。

【0005】ある端末から特定の複数の端末にパケットを送信するマルチキャスト通信の場合、ルーティングテーブルの各エントリは、送信元サブネットアドレスとアドレスマスクとマルチキャストグループアドレス、および経路情報から構成される。マルチキャスト・グループアドレスとは、送信元がパケットを送信すべき複数の宛先の集合(以下、マルチキャストグループと呼ぶ。)に対して割り当てられる識別子である。ルータは、受信パケットのヘッダ内の送信元IPアドレスと受信パケットのヘッダ内の宛先IPアドレス部に入っているマルチキャスト・グループアドレスをキーにして、ルーティングテーブルを検索する。マルチキャストの場合、検索の結果一致したエントリの経路情報は複数の送信ポート番号から構成される。ルータはこの送信ポート番号に従い、受信パケットをコピー処理して、上記特定のマルチキャストグループ向けに出力する。これらの検索処理やコピー処理の負荷が重く、マルチキャスト通信の性能が低いと、マルチキャストパケットを転送処理するルータ全体の性能低下の要因となってしまう。このため、通常のユニキャストパケット中継処理と同様に高速化が要求されている。

【0006】ルータの負荷分散処理方式による高速化に

関しては、例えば特開平6-197111号公報(以下「従来技術1」という。)に言及がある。従来技術1では、負荷分散処理方式により、高速中継可能なルータを実現することを目的としている。パケット中継を行うパケット中継モジュールがバスにより複数接続され、各々が同じくバスに接続された管理部からのルーティングテーブル情報に基づいてパケット中継機能を行う。パケット中継モジュールの増設により、性能向上ができる。さらに、従来技術1のパケット中継モジュールは、受信パケットのヘッダを抽出し、パケット転送先の検索処理を行うルーティング処理部と、受信パケットをメモリに格納し、上記解析結果に基づいて対応する他のパケット中継モジュールに受信パケットを転送する転送処理部を備え、別々の処理部で機能分担処理させることでパケット中継処理の高速化を図っている。しかし、従来技術1では、マルチキャストパケットの中継処理の高速化に関しては記載がない。

【0007】ルータにおけるマルチキャストパケットの宛先の検索(以下、マルチキャスト経路検索と呼ぶ)に関しては、例えば、(以下「従来技術2」という。)に記載されている。従来技術2では、ルーティングテーブル検索をハッシュ検索方式を用いることにより高速化している。ルーティングテーブル情報をそのまま検索すると、対象となる宛先IPアドレスが増加しテーブルのエントリが増加するに従って、検索処理にかかる検索時間が著しく増加してしまう。従来技術2は、ハッシュ値を計算するためのキーとして受信パケット内の送信元IPアドレスを使用しており、ハッシュ値が等しい宛先IPアドレスに対するエントリをグループ化したルーティングテーブル群を備える。

【0008】マルチキャストパケット受信時には、ルータ内部の経路検索処理部は、受信パケットの送信元IPアドレスからハッシュ値を計算し、該ハッシュ値に対応する上記ルーティングテーブルを検索する。ハッシュ値に従ってルーティングテーブルの検索範囲を限定することで経路検索の高速化を図っている。

【0009】一方、経路検索方式としてRadishアルゴリズムが知られている。Radishアルゴリズムに関しては、例えば“A technical memo of WIDE project, Kazuhiko Yamamoto, Akira Kato and Akira Watanabe, Radish-A Simple Table Structure for CIDR”(以下「従来技術3」という。)に記載されている。

【0010】このRadish方式のルーティングテーブルは、エントリを2分木構造に構成し、検索を高速化している。具体的には、IPアドレス形式の最上位ビットをルート(木構造の根)側とする左右にポインタを持つ複数の頂点(ノード)をポインタでつないだ2分木構造の各ノードにエントリを割り付けて構成している。この2分木構造を検索する際には、受信パケットの宛先IPアドレスを上位ビットから1ビットづつ検査し、検査されたビ

ットの値(0または1)に従って各ノードの左右のどちらかのポインタを選択して次のノードに移動する。このような検索により目的のエントリが割り付けられたノードにたどり着くことができる。

【0011】従来技術3は、宛先IPアドレスを1ビットずつ検査して2分木を辿るため、ルーティングテーブルのエントリが多くなっても、高々宛先IPアドレスのビット数回の検査を行う(2分木ノードを辿る)ことで検索を終了することができる。

【0012】上記のユニキャストの経路検索の説明で述べた、マスク長の異なる複数のエントリにマッチする場合、従来技術3のRadish方式のルーティングテーブルの検索においては、木を辿る途中で複数個のエントリにマッチすることになる。エントリがマッチする度にそのエントリ内の経路情報の候補を保持しておき、新たにエントリにマッチした場合には、以前の候補を新しい候補に更新してゆくことで最長一致検索仕様を満たすことができる。

【0013】また、従来技術3のRadish方式を適用したユニキャスト経路検索を高速に行う技術として、「IPルーティングテーブルのハードウェアによる高速検索方式」1998年電子情報通信学会通信ソサイエティ大会

(以下「従来技術4」という。)がある。従来技術4は次ノード選択回路、ネットワークアドレス比較回路、検索終了判定回路からなる。ここでネットワークアドレス比較回路は、Radish方式で構成した2分木から不要な枝を取り除いた場合、ノード自体が保持するネットワークアドレスと、受信パケットの宛先IPアドレスにノードが保持するマスクをかけたものを比較することによって、削除された枝において検査されるべきビット(1または複数の連続したビット)の値を一度に検査するための回路であり、この比較の結果が一致の場合は正しいノードに到達したことを意味し、不一致の場合は正しくないノードに到達したことを意味する。

【0014】従来技術4では、次ノードの選択と上記で説明したネットワークアドレス比較処理を並列処理し、検索処理の高速化を図っている。

【0015】また、Radish方式と同様なPatricia Trie検索アルゴリズムによるユニキャスト経路検索を高速に行う技術として、「高速IPアドレス解決H/Wエンジンの開発」1998年電子情報通信学会通信ソサイエティ大会(以下「従来技術5」という。)がある。従来技術5も従来技術4と同様にして木構造の検索の各処理を並列化して高速化を図っている。

【0016】従来技術5のPatricia Trie検索アルゴリズムによるユニキャスト経路検索をさらに高速化する技術として、「マルチウェイ化ツリーによるIPルーティングテーブルの高速検索方式」1998年電子情報通信学会通信ソサイエティ大会(従来技術6)がある。従来技術6ではPatricia Trieアルゴリズムによる2分木構造を2

のN乗分木に構成することにより、高速化を図る旨が記載されている。

【0017】また、上記のユニキャスト経路検索において最長一致検索を高速に行う技術として、特開平10-222535号公報(以下「従来技術7」という。)がある。従来技術7では、従来技術3のRadish方式とは異なる方法で2分木を構成し、その2分木を検索するハードウェアで構成されたデータ検索回路を用いて最長一致検索を高速に行っている。従来技術7の2分木の具体的な検索方法は、受信パケットの宛先IPアドレスにノード内に保持されているサブネットマスクをかけたものと、ノードのサブネットワークアドレスを比較し、比較の結果一致すればこのノードの経路情報を検索結果として検索を終了し、不一致の場合は、マスクされた宛先アドレスとノードのサブネットワークアドレスを0以上の整数とみなして大小比較を行い、その大小に従って2分木の左右の枝を辿るというものである。この際、マスク長が長いエントリから順に2分木の上部に配置しておく、その順に検索が行われるため、最長一致検索が実現できる旨が記載されている。なお、上記では、マスク長は上位何ビットがサブネットワークアドレスとして有効を示す値である、という定義に従って従来技術7の内容を説明したが、従来技術7の明細書中では、マスク長を、IPアドレスのサブネットワークアドレス部分ではない下位のビット数の長さとして定義しており、また、最長一致検索方式をベストマッチ方式と呼んでおり、上記ベストマッチ方式は、マスク長が最短のものを検索結果として採用する。

【0018】

【発明が解決しようとする課題】現在、IPマルチキャスト通信は、インターネット上に仮想的に構築されたマルチキャストネットワークにおいて限定的に利用されているにすぎない。そのネットワーク規模もそれほど大きくなく、マルチキャスト通信の検索処理や転送処理は、ルータの転送性能劣化の原因にはなっていない。しかし、今後、インターネットにおけるマルチキャスト機能を利用したアプリケーションが広く普及すると、それに伴い、ルータが保持するマルチキャストルーティングテーブルのエントリ数が多くなる。従って、ルーティングテーブルの検索処理の高速化をエントリ数の多い条件下で実現することが求められている。

【0019】従来技術2ではハッシュ方式のマルチキャスト経路検索方式が記載されているが、ハッシュ方式の場合、ハッシュ値の算出法はエントリの数にはよらず一定であり、エントリの増加に応じて柔軟にルーティングテーブル群の数を増やすことができない。従って、エントリが増加した場合、分割されるルーティングテーブル自体が長くなり、平均検索時間が長くなるという問題がある。

【0020】この点、上述のRadishアルゴリズムを用い

てユニキャストの経路検索を行う方法によれば、ハッシュ方式に比べて、ルーティングテーブルのエントリ数が増大しても、検索処理時間が増大しないという利点がある。

【0021】しかし、従来技術3乃至従来技術5には、Radishアルゴリズム(あるいはPatricia Trieアルゴリズム)を用いたマルチキャストの経路検索に関する記載はない。

【0022】また、従来技術6では、2分木構造を2のN乗分木に構成することにより、最大検索時間を1/Nに短縮する旨が記載されているが、マルチキャスト経路検索に関しては言及していない。

【0023】また、従来技術7にもマルチキャストの経路検索を高速化する方式については言及されていない。さらに、従来技術3の2分木の構造が、登録されるエントリによって一意に決定し、その最大検索時間は検索に用いるキーのビット数によって決まるのに対し、従来技術7の2分木の構造は、その作り方によっては、ノードが一方に長くつながる木ができてしまい、平均検索時間が長くなる場合がある。このような木が構成されるのを防ぐには、2分木が均等に広がるように構成するアルゴリズムが必要であるが、そのアルゴリズムについては述べられていない。また、従来技術7の2分木構造を構成する際、マスク長が長いエントリから順に2分木の上部に配置しておく、その順に検索が行われるため、最長一致検索処理が実現できると述べているが、マスク長が異なるノードのつながり方に関しては詳細は述べられていない。また、従来技術3のように検査ビットの値に従って2分木の左右の枝を辿るのではなく、ノードにおける大小比較の結果、その大小に従って2分木の左右の枝を辿るため、2のp乗分木方式を実現することができない。

【0024】そこで、本発明の第1の目的は、マルチキャストパケットの中継処理に関し、マルチキャストパケット経路検索処理の高速化の手段を提案することである。

【0025】また、本発明の第2の目的は、Radishアルゴリズムによる2分木検索方式をマルチキャスト経路検索に適用する手段を提案することである。

【0026】また、本発明の第3の目的は、検索のキーのビット数が大きい場合でも、ノードを辿る回数がビット数より少なく、最大検索時間が少ない経路検索方式を、マルチキャスト経路検索方式に適用する手段を提案することである。

【0027】本発明の第4の目的は、Radishアルゴリズムによるマルチキャストルーティングテーブルの2分木検索方式に従って経路検索処理を高速に行う回路の構成を提案することである。

【0028】

【課題を解決するための手段】 上記の目的を達成するた

めに、本発明では、以下に述べる方式および手段を持つ。

【0029】パケットの送信元の端末のIPアドレス(送信元IPアドレス)と、パケットを送信すべき端末が属するマルチキャスト・グループのグループアドレス(宛先マルチキャスト・グループアドレス)とを検索のキーとして用いるマルチキャスト経路検索において、送信元IPアドレスと宛先マルチキャスト・グループアドレスを連結して一つの経路アドレスとして定義し、この経路アドレスのビットパターンに従い、マルチキャストルーティングテーブルを2分木に構成し、受信したマルチキャストパケットの送信元IPアドレスと宛先マルチキャスト・グループアドレスとを連結した経路アドレスの上位ビットから1ビットずつ検査し、上記検査ビットの値によって2分木構造をもつ検索木を検索する。

【0030】また、上記の経路アドレスに関し、送信元IPアドレスと宛先マルチキャスト・グループアドレスの連結順を、宛先マルチキャスト・グループアドレス、送信元IPアドレスの順にして経路アドレスを定義し、この経路アドレスのビットパターンに従い、マルチキャストルーティングテーブルを2分木に構成する。

【0031】また、2分木方式の検索木を構成する各ノードを、2分木ノードから4分木、8分木、あるいは一般に2のp乗と、枝別れの数とを2のべき乗で増やすことで、マルチキャストルーティングテーブルを2のp乗分木構造に構成し、受信したマルチキャストパケットの宛先マルチキャスト・グループアドレスと送信元IPアドレスをこの順に連結した経路アドレスを、一つのノードで1ビットではなく、連続する2ビット、3ビット、あるいは一般にpビットを同時に検査し、上記連続するビットの値によって2のp乗分木構造を持つマルチキャストルーティングテーブルを検索する。

【0032】また、上記の経路アドレスとして、マルチキャスト・グループアドレス自体ではなく、マルチキャスト・グループアドレスの下位28ビットのマルチキャスト・グループIDと送信元IPアドレスをこの順に連結したものを採用し、この経路アドレスのビットパターンに従い、マルチキャストルーティングテーブルを2分木あるいは2のp乗分木に構成する。

【0033】また、検索木を構成するメモリ量を減らすために、2のp乗分木を構成する際、一つの2分木ノードと、その直下につながる(p-1)段分の合計(2のp乗-1)個分の2分木ノードを一つの2のp乗分木ノードにまとめ、まとめられる最下段の2の(p-1)乗個の2分木ノードに、それより上段のノードに割り付けられていたエントリデータを埋め込むことにより、2のp乗分木ノードを2分木ノードを2の(p-1)乗個分併せた形で構成するようにし、さらに、2分木を複数個併せるときに、共通なノードデータの要素を一つだけもつようにする。

【0034】また、2分木ノードを2の p 乗分木にまとめる際、2の p 乗分木ノードのメモリ量削減と、各2の p 乗分木ノードのメモリ量の統一のために、送信ポート情報を検索木のノードから分離して記憶手段の別領域に格納し、検索木のノードには上記送信ポート情報の格納領域へのポインタを保持することで、エントリが割り当てられている2分木ノードと、エントリが割り当てられていない2分木ノードのデータ量を同じにし、これらの2分木ノードをまとめて2の p 乗分木ノードを構成する。

【0035】また、経路アドレスの検査すべきビット位置を示すマスク長に関し、各2の p 乗分木ノードにそのノード自身のマスク長を格納するのではなく、そのノードの直下につながるノードのマスク長と格納することにより、経路検索処理を行うために検索木データが格納されている記憶手段からノードデータを読み込む際、ノードデータ全てを読み込むのではなく、直前に読込んだノードに格納されているマスク長を用いて、ノードデータのうち必要な部分のみを選択して読む。

【0036】また、各ノードの先頭に、そのノードにエントリが割り付けられているか否かを示すフラグを設け、最初に、このフラグを読み込み、エントリが割り付けられていないノードでは、経路情報を読み込まないようにする。

【0037】また、経路アドレスの上位 m ビットに対応する2の m 乗個の2の p 乗分木ノードを、記憶手段上の決まった位置に展開し、展開されたそれぞれのノードを、経路アドレスの第0ビットから第 $(m-1)$ ビットまでが取りうる値に1対1に対応させ、検索時には、経路アドレスの第0ビットから第 $(m-1)$ ビットの値に従って、該ノードを選択する。

【0038】また、宛先アドレスと送信元アドレスをこの順に連結して経路アドレスとし、上記経路アドレスの上位 m ビットに対応する2の m 乗個の2の p 乗分木ノードを決まった位置に格納し、それ以降の木は2の p 乗分木構造に構成して格納する記憶手段を持ち、経路検索時には、受信したマルチキャストパケットの宛先アドレスと送信元アドレスをこの順に連結した経路アドレスの第0ビットから第 $(m-1)$ ビットの値に従って、メモリ上の決まった位置に展開されたノードの一つを選択し、経路アドレスの第 m ビット以降は経路アドレスを p ビットずつ検査して2の p 乗分木に構成されたマルチキャストルーティングテーブルを検索する回路をもつ。

【0039】その他の本願が解決しようとする課題、及びその解決手段は、後述の「発明の実施の形態」の欄、及び本願図面で明らかにされる。

【0040】

【発明の実施の形態】以下、本発明の実施例を図面を用いて説明する。

【0041】まず、一般的なネットワークの構成および

そのネットワークにおけるパケット中継動作を図2、3、4を用いて説明する。

【0042】図2に示したネットワークは、パケットを送受信する端末(T11、T12、T21、T22、T31、T32、T41、T42、T51、T52)と、端末を相互に接続するサブネットワーク(SN1～SN5)と、該サブネットワークを相互に接続するルータ(R1～R3)から構成される。また、ルータR1はポートP11、P12、P13を持ち、ルータR2はポートP21、P22、P23、P24を持ち、ルータR3はポートP31、P32、P34、P35を持つ。各端末には識別子としてアドレスが割り当てられており、例えばIPネットワーク上の各端末には32ビットのIPアドレスが割り当てられている。このIPアドレスについて図3(a)を用いて説明する。図3(a)に示したIPアドレス(1)は、サブネットワークアドレス(2)とホストアドレス(3)から構成される。

【0043】サブネットワークアドレスは、サブネットワークの識別子であり、ホストアドレスはサブネットワーク内の端末の識別子である。サブネットワークアドレスのビット数はサブネットマスク長で示されており、このサブネットマスク長だけ上位ビットから有効なマスク(サブネットマスク)とIPアドレスのビットごとの論理積をとることにより、IPアドレスからサブネットワークアドレスを得ることができる。

【0044】次に、ある端末から特定の一つの端末にパケットを送信するユニキャスト通信について、図2を用いてパケット中継動作を説明する。一例として、端末T11から端末T22へユニキャストパケットを送信する場合について説明する。送信されるパケットのヘッダフォーマットを図3(c)に示す。図3(c)の送信元IPアドレスフィールド(7)には送信元端末T11のIPアドレスが格納され、宛先IPアドレスフィールド(8)には宛先端末T22のIPアドレスが格納される。端末T11からのユニキャストパケットを受信したルータR1は、パケットヘッダ内の宛先IPアドレスフィールドを検査し、宛先端末T22がサブネットワークSN2上にあり、かつ、サブネットワークSN2がルータR2経由で接続されていることを認識し、パケットを送信すべきルータR2のIPアドレス(ネクストホップIPアドレスと呼ぶ)および送信すべきポートP12を決定し、パケットを出力する。ルータR1からのユニキャストパケットを受信したルータR2は、パケットヘッダ内の宛先IPアドレスフィールドを検査し、宛先端末T22がサブネットワークSN2上にあり、かつ、サブネットワークSN2がルータR2に直接接続されていることを認識し、ネクストホップIPアドレス(今の場合、端末T22のIPアドレス)および送信すべきポートP22を決定し、パケットを出力する。

【0045】次に、ルータの構成およびルータのユニキ

キャストパケットの中継動作を図4を用いて説明する。図4のルータ10は、入力ポート11-i(=1~N)、出力ポート12-i(=1~N)、ルーティング処理部13-i(=1~N)、スイッチ14、通信制御部15-i(=1~N)、管理部16から構成される。上記で説明したユニキャストパケットが入力ポート11-iから入力されると、ルーティング処理部13-iがパケットのヘッダ内の宛先IPアドレスをキーにしてユニキャストルーティングテーブルを検索する。この検索により、ネクストホップIPアドレスと、出力すべき出力ポート12-iの番号(送信ポート番号)を決定し、上記ネクストホップIPアドレスと送信ポート番号をパケットに付加してスイッチ14に送出する。スイッチ14は上記送信ポート番号に対応する出力ポート12-iを保持する通信制御部15-iにパケットをスイッチングする。スイッチングされたパケットを受信した通信制御部15-iは、パケットに付加されているネクストホップIPアドレスに対応するデータリンク層のMAC(Media Access Control)アドレスを対応テーブルに従って決定してパケットの宛先MACアドレスとして付加して、送信ポート番号に対応する出力ポート12-iにパケットを送出する。

【0046】なお、上記のユニキャストルーティングテーブルは、管理部16が他のルータと接続情報をやり取りして作成し、各ルーティング処理部13-iに配布する。また、管理部16は、他のルータおよび端末とデータリンク層の情報のやり取りをし、ネクストホップIPアドレスと、そのIPアドレスを持つルータおよび端末のMACアドレスの対応テーブルを作成して各通信制御部15-iに配布する。

【0047】次に、ある端末から特定の複数の端末にパケットを送信するマルチキャスト通信について、図2を用いてパケット中継動作を説明する。一例として、端末T11からある特定のマルチキャストグループに属する端末(T21、T31、T32、T41、T42、T51)へマルチキャストパケットを送信する場合について説明する。マルチキャストグループは、クラスDのIPアドレス(マルチキャスト・グループアドレス)によって識別される。このクラスDのIPアドレスの構造を図3(b)に示す。クラスDのIPアドレス(4)は、クラスDであることを示す"1110"という上位4ビット(5)と28ビットのマルチキャスト・グループID(以下、GID)(6)から構成される。上記のマルチキャスト・グループアドレスによって識別されるマルチキャストグループは、複数のサブネットワークにまたがる可能性がある。また、グループに属する端末のメンバー構成は動的で、端末は自由にマルチキャストグループに参加、離脱することができる。

【0048】端末T11から、あるマルチキャストグループ宛に送信されたマルチキャストパケットのヘッダ内の送信元IPアドレスフィールド(図3(c)の7)には送

信元端末T11のIPアドレスが格納され、宛先IPアドレスフィールド(図3(c)の8)には、ユニキャストパケットの場合とは異なり、宛先のマルチキャスト・グループアドレスが格納される。端末T11からのマルチキャストパケットを受信したルータR1は、ユニキャストパケットの場合とは異なり、パケットヘッダ内の送信元IPアドレスフィールドと宛先IPアドレスフィールドの2つのフィールドを検査し、このマルチキャストパケットの送信元端末T11がサブネットワークSN1上にあることを認識し、さらに、このサブネットワークSN1から、宛先マルチキャストアドレスに属する端末の全てにパケットを送信するためにルータR1がパケットを送信すべき送信ポートP12、P13を決定し、上記ポートにパケットをコピーして出力する。ルータR1からのマルチキャストパケットを受信したルータR2は、R1と同様な検査を行い、このマルチキャストパケットの送信元端末T11がサブネットワークSN1上にあることを認識し、さらに、宛先マルチキャストアドレスに属する端末が存在するサブネットワークSN2、SN3にパケットを送信するためにルータR2がパケットを送信すべき送信ポートP22、P23を決定し、上記ポートにパケットをコピーして出力する。ルータR3も同様にR1から受信したパケットを送信すべきポートP34、P35を決定し、上記ポートにパケットをコピーして出力する。

【0049】次に、ルータのマルチキャストパケットの中継動作を図4を用いて説明する。上記で説明したマルチキャストパケットが入力ポート11-iから入力されると、ルーティング処理部13-iがパケットのヘッダ内の送信元IPアドレスと宛先IPアドレス(宛先マルチキャストアドレス)をキーにしてマルチキャストルーティングテーブルを検索する。この検索により、出力すべき1つまたは複数の出力ポート12-iの番号(送信ポート番号)を決定し、受信したパケットをコピーし、各送信ポート番号をパケットに付加してスイッチ14に送出する。スイッチ14は上記送信ポート番号に対応する出力ポート12-iを保持する通信制御部15-iにパケットをスイッチングする。スイッチングされたパケットを受信した通信制御部15-iは、ユニキャストパケットの場合とは異なり、マルチキャストパケットのヘッダ内の宛先マルチキャストアドレスを用いてデータリンク層のマルチキャスト用MACアドレスを生成し、パケットの宛先MACアドレスとして付加して、送信ポート番号に対応する出力ポート12-iにパケットを送出する。

【0050】なお、上記のマルチキャストルーティングテーブルは、ユニキャストルーティングテーブルと同様に管理部16が他のルータと接続情報をやり取りして作成し、各ルーティング処理部13-iに配布する。

【0051】次に、本発明の一実施例として、ルータの全体構成例を図1を用いて説明する。

【0052】まず、ルータ700の構成を図1を用いて説明する。505はスイッチであり、図4のスイッチ14に対応する。上記スイッチ505に、装置全体の管理機能と経路情報の収集・分配等の機能をもつ管理部601が接続されている。管理部601は図4の管理部16に対応する。また、スイッチ505には、高速にパケット中継処理を行う機能を有するパケット処理部600を複数接続することができる。パケット処理部600は、図4のルーティング処理部13-iに対応する。管理部601は、それぞれのパケット処理部600にルーティングテーブルを配布し、それぞれのパケット処理部600が上記ルーティングテーブルに基づいてパケット中継処理を行うことによりルータの機能を実現する。

【0053】さらに、各パケット処理部600には、ネットワークを接続するための通信ポートを備える通信制御部501、602、603(図4の通信制御部15-iに対応)が接続される。通信制御部は、各種のネットワークに接続することができる。また、そのネットワークの通信速度に応じて、高速通信(例えば、図2の500)ならば1ポート、低速通信(例えば、図2の604、605)ならば複数ポートを設け、1つまたは複数のネットワークと接続することができる。

【0054】次に、上記のパケット処理部600の構成について、図1を用いて説明する。図1において、パケット処理部600は、パケット転送処理を行う転送処理部504と、パケットを格納するパケットバッファを備えるメモリ503と、メモリ503に格納されたパケットのヘッダ情報を用いて経路検索を行う経路検索処理部400と、転送処理部504および通信制御部501を制御する制御回路502とを有する。経路検索処理部400は、ルーティングテーブルを格納するメモリ300と、そのルーティングテーブルを検索し、検索結果を転送処理部504へ出力する経路検索処理回路200から構成される。

【0055】次に、パケット処理部600における、マルチキャストパケットの中継処理の動作を説明する。なお、パケット中継処理を行う前に、管理部601が、スイッチ505に接続している全パケット処理部600に、それぞれのマルチキャストルーティングテーブルを配布しており、各パケット処理部600はメモリ300にルーティングテーブルを格納している状態とする。通信制御部501がネットワーク500から受信したマルチキャストパケットをメモリ503へ転送し格納する。転送処理部504はメモリ503に格納されたパケットデータの内、パケットヘッダ情報を抽出して経路検索処理部400へ転送する。経路検索処理回路200は、受信パケットのヘッダ内の送信元アドレスと宛先マルチキャスト・グループアドレスを用いてメモリ300内に格納されているマルチキャストルーティングテーブルの検索を行い、検索結果として、パケットを送信すべきN個

(Nは1以上の整数)の送信ポート番号を転送処理部504へ転送する。

【0056】上記の検索結果を受信した転送処理部504は、先にメモリ503に格納されているパケットに対してパケットヘッダの書換などの処理を施し、パケットに上記N個の送信ポートの内の1つの番号を付加してスイッチ505に転送する。転送処理部504は上記の処理をN回繰り返し、1つの受信パケットに対する転送処理を終了する。また、送信ポート番号が自パケット処理部に接続されている通信制御部のポート番号である場合は、自パケット処理部に接続されている通信制御部にパケットを送信する。

【0057】スイッチ505は転送処理部504により転送されたパケットを、付加されている送信ポート番号に従い、各パケット処理部にスイッチングする。上記スイッチングされたパケットを受信したパケット処理部600は、上記パケットを通信制御部(501、602、603)へ転送する。通信制御部は、上記パケットのヘッダ内の宛先マルチキャストアドレスからマルチキャスト用MACアドレスを生成し、宛先MACアドレスとしてパケットに付加して各ネットワークに送信する。

【0058】本発明は、上記の経路検索処理部400の高速化に関するものである。以下では、この経路検索処理部400に関し、まず、方式1として、マルチキャストルーティングテーブルを2分木に構成し、その2分木構造を検索する検索方式(以下、2分木検索方式と呼ぶ)について説明する。方式2として、方式1を改良し、さらに検索を高速に行うために、マルチキャストルーティングテーブルを2のp乗分木に構成し、その2のp乗分木構造を検索する検索方式(以下、2のp乗分木検索方式と呼ぶ)について説明する。また、その方式2の説明の中で、マスク長mビットのノードをメモリ上に展開して方式1、2のノードを渡る回数を削減する方式(方式3)について説明する。最後に、本発明の一実施例として、方式2および3の方式を適用した経路検索処理を実現するハードウェア構成について説明する。

【0059】まず、方式1の2分木検索方式について説明する。

【0060】従来の技術で説明したように、マルチキャスト経路検索を行う際、受信パケットのヘッダ内の送信元IPアドレスと、宛先マルチキャスト・グループアドレスを検索のキーとしてマルチキャストルーティングテーブルを検索する。すなわち、この2つの情報の組み合わせにより、そのマルチキャストパケットの送信元から宛先までの経路が識別される。

【0061】従って、送信元IPアドレスとマルチキャスト・グループアドレスを連結して、これを経路アドレスとして定義し、この経路アドレスと経路情報の組をエントリとしてマルチキャストルーティングテーブルを構成することにより、この経路アドレスをキーとして、マ

ルチキャストルーティングテーブルを検索することができる。

【0062】また、上記の2つのアドレスのうち、マルチキャスト・グループアドレスに関しては、マルチキャスト・グループ一つに対して一つのアドレスが割り当てられているので、ルーティングテーブル検索時は、エントリに保持されたマルチキャスト・グループアドレスとの一致比較を行えばよい。

【0063】しかし、上記の2つのアドレスのうち、送信元IPアドレスに関しては、受信パケットの送信元IPアドレスに、エントリに保持されているサブネットマスク長だけ上位ビットから有効なマスクをかけたものとエントリに保持されている送信元サブネットワークアドレスに一致するエントリの経路情報を検索結果とするが、この場合、ユニキャスト経路検索と同様に最長一致検索を行う必要がある。

【0064】以上のことから、送信元IPアドレスとマルチキャスト・グループアドレスの連結順を、マルチキャスト・グループアドレス、送信元IPアドレスの順にして、これを経路アドレスとして定義し、この経路アドレスのビットパターンに従ってマルチキャストルーティングテーブルを2分木構造に構成することにより、従来技術3に述べられているRadish方式を用いたユニキャスト経路検索方式と同様にして、経路アドレスに関して最長一致検索を行うことができる。この場合、経路アドレスの定義から、マルチキャスト・グループアドレスについては一致比較、送信元IPアドレスに関しては最長一致検索を行うことができる。

【0065】図5に2分木で構成されたマルチキャストルーティングテーブルを検索する際に用いる経路アドレスを示す。なお、この経路アドレスは後述する2のp乗分木検索方式にも採用する。図5(a)の経路アドレス32は、マルチキャスト・グループアドレス(30)と送信元IPアドレス(31)とを連結したものであり、ビット数は64ビットである。図5(b)の経路アドレス33は、マルチキャスト・グループアドレス(30)から"1110"というIPアドレスのクラスDを示す上位4ビットの値を除いたマルチキャスト・グループID(34)と、送信元IPアドレス(31)とを連結したものであり、ビット数は60ビットである。木構造を構成するための値(今の場合は検索アドレス)のビット数が少ない方が木の構造が浅くなり、検索時に辿る最大ノード数が減るため、以下では、ビット数の少ない図5(b)の経路アドレス33を採用して説明する。

【0066】次に、上記の経路アドレスを用いてマルチキャストルーティングテーブルを2分木に構成する方法について説明する。木構造の考え方は経路アドレスのビット数には依存しないので、簡単のため、マルチキャスト・グループIDは1ビット、送信元IPアドレスは2ビットと仮定して説明する。この場合経路アドレスは3

ビットとなる。

【0067】図6に2分木の構造の例を示す。図6に示すように、各ノードは、経路アドレス $t a a t$ と、経路アドレスマスク長 n 、および、そのノードの下につく二つのノードへのポインタを持つ。ここで、経路アドレスマスク長とは、そのノードが持つ経路アドレスの上位何ビットが有効であることを示す値であり、また、そのノードでの受信パケットの経路アドレスの検査ビット位置も示す。また、経路アドレスと経路アドレスマスクの組を表記法 $t a a : n t$ で表す。この表記法により、マスク付き経路アドレスを表すことにする。

【0068】各ノードを木の上から順に、経路アドレスマスク長0ビット、1ビット、2ビット、3ビットのノードと呼ぶ。

【0069】経路アドレスマスク長0ビットのノード(40)では、経路アドレスマスク長の値で示されている経路アドレスの第0ビットを検査し、その値が0か1かに従い左/右のポインタを辿ることにより経路アドレスマスク長1ビットのノード(41、42)に移り、経路アドレスマスク長1ビットのノードでは、経路アドレスの第1ビットが0か1かに従い左/右のポインタを辿ることにより経路アドレスマスク長2ビットのノード(43、44、45、46)に移り、経路アドレスマスク長2ビットのノードでは、経路アドレスの第2ビットが0か1かに従い左/右のポインタを辿ることにより経路アドレスマスク長3ビットのノード(47、48、49、50、51、52、53、54)に移る。

【0070】検索したい経路アドレスについて、この木の経路アドレスマスク長0ビットのノード(40)から順に、各ビットが0か1かに従いポインタを辿った場合、経路アドレスマスク長0ビットのノードは経路アドレスがどの値をとる場合にも通過し、経路アドレスマスク長1ビットのノード(41、42)は左から順に経路アドレスの各ビットが0xx、1xxの場合に通過し、経路アドレスマスク長2ビットのノード(43、44、45、46)は左から順に経路アドレスの各ビットが00x、01x、10x、11xの場合に通過し、経路アドレスマスク長3ビットのノード(47、48、49、50、51、52、53、54)は左から順に経路アドレスの各ビットが000、001、010、011、100、101、110、111の場合に通過する。ここで、xは、そのビット値が0または1のどちらでも良いことを示す。

【0071】従って、経路アドレスマスク長0ビットのノード(40)は、経路アドレスがマスク付き経路アドレス000:0に属する場合に通過し、経路アドレスマスク長1ビットのノード(41、42)は、経路アドレスがそれぞれマスク付き経路アドレス000:1、100:1に属する場合に通過し、経路アドレスマスク長2ビットのノード(43、44、45、46)は、経路

アドレスがそれぞれマスク付き経路アドレス000:2、010:2、100:2、110:2に属する場合に通過し、経路アドレスマスク長3ビットのノード(47、48、49、50、51、52、53、54)は、経路アドレスがそれぞれマスク付き経路アドレス000:3、001:3、...、111:3に属する場合に通過する。

【0072】上記の通り、この木の各ノードは、経路アドレスと経路アドレスマスク長が異なる全マスク付き経路アドレスに1対1に対応している。

【0073】上記の2分木において、図7に示すマルチキャストルーティングテーブルの各エントリに対応するノード40、48、44、42、53に該エントリを割り付ける。具体的には、図12を用いて後述するように、ノードのデータ構造内にエントリが割り付けられているか否かを示すフラグを設け、エントリが割り付けられたノード内の上記フラグの値を1に設定し、かつ送信ポート情報あるいは送信ポート情報へのポインタをノード内に格納する。本明細書では、ノードに*を付記することにより、ノードにエントリが割り付けられていることを表すことにする。

【0074】受信パケットの経路アドレスが011(図6の60)の場合、この木の上から各ビットが0か1かに従いポインタを辿ると、*が付記されたノード40、44が、マスク付きの検索で一致することが分かる。そこで、ルーティングテーブルエントリが複数一致した場合は、送信元IPアドレスに関して最長一致検索を行うことを考慮すると、一致した*付きノード40、44の内、上記で定義した経路アドレスマスク長が最も長いノード、すなわち最も末端に近いノード44に割り付けられたエントリの経路情報を、ルーティングテーブルの検索結果とする。なお、図7のルーティングテーブルの例では、エントリ内の経路情報は省略している。

【0075】上記検索方法から分かるように、*が付記されておらず、かつ*付きのノードにたどり着くための途中経路にもなっていないノード47、49、50、51、52、54、45は木から取り除いても、検索結果には影響しない。むしろ、最下のノードが、*が付記されていないノードの場合は最下まで移動せずに検索が終了するために効率的である。そこで、*が付記されておらず、かつ*付きのノードにたどり着くための途中経路にもなっていないノードを木から取り除くと図8のようになる。

【0076】上記の方法を用い、経路アドレス長が32ビットと仮定した場合、図9に示す経路テーブルに対応する2分木を構成すると、図10に示す分岐も*も無いノードの長い列61ができる。なお、図9のルーティングテーブルの例では、図7と同様にエントリ内の経路情報は省略している。また、図10の0xNNNNNNNという

表記法は、NNNNNNNが16進数表示であることを示す。このように、左右の片方のポインタだけに次のノードが繋がりを、かつエントリが割り当てられていないノードを取り除くことにより、検索時に辿るノード数を減らすことができる。

【0077】図10の例においては、分岐も*も無いノード列61を取り除き、直ぐ上のノード62の分岐方向(図10では右側)に、取り除かれたノード列61の直ぐ下のノード63を付ける。その結果、図11に示す形となる。このように途中のノード列を取り除くことを、以後、木の縮退と呼ぶ。

【0078】次に縮退した2分木の検索法を説明する。

【0079】図11に示す例では、経路アドレスマスク長0ビットのノード62で第0ビットの検査を行った後、経路アドレスマスク長15ビットのノード63に跳ぶので、経路アドレスマスク長15ビットのノード63で第15ビットだけを検査したのでは、途中のビット、即ち第1ビットから第14ビットが検査できない。そこで、第1ビットから第15ビットの検査を一回の処理で行う為に、受信パケットの経路アドレスの第1ビットから第15ビットとノード63の経路アドレス0x85040000の第1ビットから第15ビットの一致比較を行う。比較結果が一致すれば正しいノードにたどり着いたこと、即ち、縮退しない木で1ビットずつ比較しても、このノードにたどり着いたことを意味し、一致しなければ、正しくないノードにたどり着いたこと、即ち縮退しない木では、行き先が無いことを意味する。

【0080】ここで、図11に示す例では、第0ビットは、既にテストされ、第0ビットが等しくなる方の分岐が選択されている為、常に一致する。一般に、あるノードにたどり着く毎に、正しいノードにたどり着いたか否かを検査していれば、第0ビットから、そのノードの経路アドレスマスク長までのビットに関しては、受信パケットの経路アドレスと、ノードの経路アドレスとが等しいことが保証されているので、次のノードにたどり着いたときに、前にどのビットまで検査したかに関らず、第0ビットからノードの経路アドレスマスク長までのビットに関して、受信パケットの経路アドレスとノードの経路アドレスとが等しいか否かを調べて良い。

【0081】次に、上記で説明した2分木を構成するノードのデータ構造について、図12を用いて説明する。図12(a)において、ワード110、111内の次のノードのマスク長0、1は、このノード自身の経路アドレスマスク長ではなく、このノードの直下のノードの経路アドレスマスク長である。このように、自分自身でなく直下のノードのマスク長を設定する理由は高速化のためであり、その説明は図20を用いて後述する。ワード110、111内のフラグ0、フラグ1は、このノードがエントリが割り付けられているノードか否か、即ち、図6、8、10、11で示した木の例において、このノ

ドが、"*"が付記されているノードか否かを示すフラグ(以下、エントリ有りフラグと呼ぶ)や、このノードに一致するパケットに関する属性フラグなどが設定されるフィールドである。フラグ0とフラグ1には、同じ値を設定する。これは、ワード110とワード111の一つだけを読めば良いようにするためであり、このように、ノードの全てを読むのではなく、一部分を読むことによる高速化については図20を用いて後述する。ワード110、111内の次のノードへのポインタ0、1は、経路アドレスの、このノードの経路アドレスマスク長で示されるビット位置の値が、それぞれ0、1のときに次に辿るノードへのポインタである。ワード112、113の経路アドレスは、このノードに対応する経路アドレスであり、例えばワード112にはノードの経路アドレスの内、上位28ビットを設定し、ワード113にはノードの経路アドレスの内、下位32ビットを設定する。

【0082】図12(a)では、2分木ノードのデータ構造の中に送信ポート情報114を保持する例を示している。図12(b)では、送信ポート情報114を2分木ノード構造から分離し、分離した送信ポート情報へのポインタ115を保持する例を示している。送信ポート情報は、受信マルチキャストパケットを転送すべき1つまたは複数の送信ポート番号を示し、ルータのポート数が多い場合は数ワード必要になる。従って、図12(a)のように、2分木ノード構造内に送信ポート情報114を保持する場合、エントリが割り付けられている2分木ノードと、エントリが割り付けられていない2分木ノードで、その大きさが大きく異なることになる。このようにノードの大きさが異なると、ノード内の次ノードへのポインタに必要なビット数が多くなることや、そのポインタから実際のメモリのアドレスを計算する処理が複雑になり、特にハードウェア構成が複雑になるという問題が考えられる。

【0083】図12(b)のように送信ポート情報をノードから分離し、ノードにはその送信ポート情報へのポインタ115を保持する場合、エントリが割り付けられている／いないに関わらずノードの大きさは一定になり、ハードウェア構成を単純にする利点がある。

【0084】次に、方式2の2のp乗分木検索方式について説明する。

【0085】上記で説明した2分木方式では、一つのノードに2つの分岐先があり1ビットずつ検索していた。検索のキーのビット数は、検索時に辿るノードの数(の最大値)に一致する。したがって、最大の検索時間は(式1)で表される。

【0086】(最大検索時間)=(一つのノードの処理時間)×(検索のキーのビット数) …(式1)
すなわち、最大検索時間は検索のキーのビット数に比例する。この最大検索時間をさらに大幅に短縮し、高速化することが求められる。

【0087】そこで、以下で説明する2のp乗分木方式では、一つのノードに2のp乗の分岐先を設け、同時にpビットずつ検索することにより、従来に比べ、検索時間を1/pに短縮するものである。以後、一つのノードに2のp乗の分岐先があるノードのことを2のp乗分木ノードと呼ぶ。

【0088】2のp乗分木ノードは、2分木方式で説明した2分木ノードから構成される木を変形することにより作成する。木の変形の方法は、経路アドレスマスク長nの2分木ノード一つと、この2分木ノードの下に存在する、経路アドレスマスク長(n+1)ビットから(n+p-1)ビットの2分木ノードを、一つの2のp乗分木ノードに対応させるものである。また、以下では、上記のように対応させてまとめた2のp乗分木ノードに関し、まとめる前に頂点にあった2分木ノードのマスク長nを引き継いで、この2のp乗分木ノードのマスク長をnと定義する。

【0089】上記の変形法の例として、2分木から8分木への変形法を図13に示す。

【0090】経路アドレスのビット数が60ビットの場合で、8分木を構成する場合を考える。この場合、一つの8分木に対応させる2分木ノードの選び方として、マスク長0~2、3~5、……、57~59、60ビットの2分木ノードを、それぞれ一つの8分木ノードとする場合(図13(a))と、マスク長1~3、4~6、……、55~57、58~60ビットの2分木ノードを、それぞれ一つの8分木ノードとする場合(図13(b))と、第2~4、5~7、……、56~58、59~60ビットの2分木ノードを、それぞれ一つの8分木ノードとする場合(図13(c))の3通りがあり、どの区切り方でも、構成可能だが、木全体では、エントリの追加、削除を容易に行えるように、上記3通りのビット位置の区切り方の内、一つを使用する。

【0091】上記3通りのビット位置の区切り方の内、最初の区切り方以外では、マスク長が0ビットから始まっていないので最初のビットの検査を別に行う必要がある。この検査には、マスク長mビットのノードをメモリ上の決まった位置に展開する方法(以下、方式3と呼ぶ)を使用する。

【0092】以下、方式3について説明する。方式3は、2のp乗分木構造のマスク長mビットのノードを、ノードがある場合も無い場合も区別せず全て予め用意し、メモリ上の決まった位置に展開しておき、経路検索開始時に、受信パケットの経路アドレスの上位mビットの値に従って上記の展開されたノードの内の一つを選択してメモリから読み込む。この方式により上位mビット分の検索時間を省くことができる。

【0093】上記の方式3を適用した例を、図13(b)、(c)に示した区切り位置の場合について図14(a)、(b)に示す。

【0094】図14(a)に示す構成では、マスク長1ビットの8分木ノード90、91をメモリ上の決まった位置に並べ、それぞれを、受信パケットの経路アドレスの第0ビットが0か1かに従い選択することにより、最初の第0ビットの検査を行ったことになる。図14(b)に示す構成では、マスク長2ビットの8分木ノード92、93、94、95をメモリ上の決まった位置に並べ、それぞれを、受信パケットの経路アドレスの第0～1ビットの値が00か01か10か11かに従い選択することにより、最初の第0、1ビットの検査を行ったことになる。

【0095】さらに、図13(a)、(b)、(c)に示す各ビットの区切り位置の場合に、それぞれ最初に並べる8分木ノード数を1、2、4個ではなく、これらの8倍である8、16、32個、8の2乗倍である64、128、256個、或いは一般に8の q 乗倍個にし、最初の1回、2回、或いは一般に q 回の8分木ノードの検索時間を無くすことも可能である。この場合、メモリ上の決まった位置に展開する8分木ノードのマスク長 m は、図13(a)、(b)、(c)に示す各ビットの区切り位置の場合、それぞれ、 $m=3 \times q$ 、 $1+3 \times q$ 、 $2+3 \times q$ となり、この上位 m ビット分の検索時間を省くことができる。 $q=1$ の場合、即ち8分木ノードの1回のノードの検索時間を無くす場合で、図13(a)、(b)、(c)に示す3通りの各ビットの区切り位置の場合についてのメモリ上へのノードの展開法を、図15(a)、(b)、(c)にそれぞれ示す。

【0096】図15(a)では、マスク長 $m=3$ の8分木ノードを8個メモリ上の決まった位置に展開し、受信パケットの経路アドレスの第0～2ビットの計3ビットの値に従って、展開された8分木ノードの内の一つを選択する。図15(b)では、マスク長 $m=1+3=4$ の8分木ノードを16個メモリ上の決まった位置に展開し、受信パケットの経路アドレスの第0～3ビットの計4ビットの値に従って、展開された8分木ノードの内の一つを選択する。図15(c)では、マスク長 $m=2+3=5$ の8分木ノードを32個メモリ上の決まった位置に展開し、受信パケットの経路アドレスの第0～4ビットの計5ビットの値に従って、展開された8分木ノードの内の一つを選択する。

【0097】以上、8分木を例にして方式3について説明したが、同様にして、マスク長 m ビットの2の p 乗分木ノードをメモリ上の決まった位置に展開し、上位 m ビット分の検索時間を省くことができる。 p 、 q の値を大きくすると、経路検索時間を短縮することができるが、多くのメモリを必要とするので、 p 、 q の値は、メモリ効率と性能のトレードオフから決めるようにする。

【0098】以上、方式3について、2の p 乗分木ノードの場合について述べたが、同様にして、方式1の2分木検索方式にこの方式3を適用することも可能である。

【0099】次に、4分木ノード、8分木ノード、16分木ノード、或いは一般に2のべき乗分木のノードの構成法を図16を用いて説明する。

【0100】図16は4分木で、ある1つの2分木ノードA、B、C、D、Eと、その直下の2個の2分木ノードA0、A1、B0、B1、C0、C1、D0、D1、E0、E1の、各々合計3個の2分木ノードをまとめて一つの4分木ノード100、101、102、103、104にする場合の例であり、合計3個の2分木ノードをつぶして、下の方の2分木ノードだけの大きさにする。つぶし方は、経路アドレスマスク長に関する最長一致検索の仕様に従い、経路検索を行った場合に、2分木の場合と4分木の場合とで、経路検索結果が同じになる、という条件を満たすように行う。

【0101】4分木の場合の、このノードのつぶし方を図17に示す。図17において、*Aという表記は、ノードAに割り付けられているエントリ内の経路情報を表すことにする。また、*Aという表記がないノードにはエントリが割り付けられていないことを意味する。

【0102】2分木構造において、3つのノードが全部ある場合(図17(a))、全ノードにエントリが割り付けられていたら、下のノードの経路情報*A0、*A1を残し、上のノードは削除する。これは、ノードAのエントリが一致したらノードA0かノードA1のどちらかのエントリが必ず一致するので、最長一致検索を行うことから、ノードAの経路情報*Aが使用されることが無いからである。

【0103】上のノードAにエントリが割り付けられており、下のノードA0、A1の内A1にだけエントリが割り付けられていない場合(図17(b))、ノードA1に、Aの経路情報*Aを格納する。下のノードA0、A1の内A0だけエントリが割り付けられていない場合も同様である。

【0104】下のノードA0、A1の両方にエントリが割り付けられていない場合(図17(c))には、A0、A1の両方のノード内に、Aの経路情報*Aを格納する。

【0105】下のノードA1が無い場合(図17(d))、A1を補い、A1には、上のノードAの経路情報*Aを格納し、さらに、ノードA1の下にはノードが繋がっていないので、ノードA1の下の方のポインタはNULLにする。下のノードA0、A1の内A0だけエントリが割り付けられていない場合も同様である。

【0106】下のノードA0、A1の両方が無い場合(図17(e))、両方を補い、両方のノードに、Aの経路情報*Aを格納し、両方のノードの下の方のポインタはNULLにする。

【0107】上のノードAにエントリが割り付けられていない場合(図17(f))、ノードAをただ単に削除する。

【0108】上のノードAに経路が割り付けられておら

ず、下のノードA0、A1の内、A1にもエントリが割り付けられていない場合(図17(g))、4分木にした場合もA1の経路情報は無い。下のノードA0、A1の内、A0にエントリが割り付けられていない場合も同様である。

【0109】上のノードAにエントリが割り付けられてなく、下のノードA0、A1の両方にエントリが割り付けられていない場合(図17(h))、4分木にした場合もA0、A1の両方のノードの経路情報は無い。

【0110】下のノードA0だけしかない場合(図17(i))には、下のノードA1を補う。下のノードA1だけしかない場合も同様である。

【0111】8分木の場合も同様にして、一つにまとめる7個の2分木ノードをつぶして、一番下の4個のノードだけの大きさにする。上の方の3つの2分木ノードのつぶし方の2つの例を図18に示す。

【0112】図18(a)は一つにまとめる7個の2分木ノードが全てあるが、その内のいくつかにはしかエントリが割り付けられていない場合の例であり、最下の4つのノードの内、エントリが割り付けられていないノードA01、A10には、そのノードの上方につながっているノードの内、エントリが割り付けられている最も下、即ち最も経路アドレスマスク長が長いノード(それぞれ、A、A1)の経路情報*A、*A1を格納する。

【0113】図18(b)は一つにまとめる7個の2分木ノードの内のいくつかしかノードが存在しない場合の例であり、存在しないノードA01、A10を、まずエントリが割り付けられていないノードとして補い、図18(a)と同じ規則で、経路情報を格納する。最下の4つのノードA00、A01、A10、A11の内、補ったノードA01、A10の下には、ノードが繋がっていないので、下のノードへのポインタはNULLにする。

【0114】一般に2のp乗分木の場合も同様にして、一つにまとめる(2のp乗-1)個の2分木ノードをつぶして、一番下の2の(p-1)乗個のノードだけの大きさにする。

【0115】また、4、8、16、...分木ノードでは、2分木ノードを2、4、8、...個まとめて扱うので、1つの4、8、16、...分木ノードにまとめられる2分木ノード間で共通な要素は、一つだけ持てばよく、これにより、4、8、16、...分木ノードのメモリ量を小さくできる。1つの4、8、16、...分木ノードにまとめられる2分木ノード間で共通な要素には、ノードが持つ経路アドレス、及び、経路アドレスマスク長があるが、経路アドレスマスク長については、後述するように、このノード自身の経路アドレスマスク長ではなく、このノードの直下のノードの経路アドレスマスク長を持つようにするので、メモリ量を小さくする効果は無い。

【0116】次に、2のp乗分木ノードの一つの例とし

て、4分木ノードのデータ構造について図19を用いて説明する。4分木ノードを構成する際、上記2分木ノードを2つまとめ、まとめられる2つの2分木ノードに共通のデータの一つだけ保持するようにする。この共通のデータは、そのノードの経路アドレスだけである(ワード124、125)。また、まとめられる2分木ノードが、経路が割り付けられているか否かによりその大きさが異なると、それらをまとめて構成する4分木ノードの大きさも異なってしまう。従って、図12(b)で示したように、送信ポート情報をノードから分離し、ノード内にはこの送信ポート情報へのポインタを保持する2分木ノードをまとめる。

【0117】図19において、ワード120、121、122、123内の次ノードのマスク長tt、フラグtt、次ノードへのポインタttはそれぞれ、受信バケットの経路アドレスの、このノードの経路アドレスマスク長mで示された第mビット、および第m+1ビットの値ttに対応する値である。ここでttは00、01、10、11の値を取りうる。2分木の場合と同様に、フラグ00とフラグ01は同じ値を設定し、フラグ10とフラグ11も同じ値を設定する。

【0118】図12(b)に示す例では、2分木ノードは、2のべき乗の大きさである16バイトにわずかに入りきらない大きさになっているが、4分木ノードにし、1ノード内に経路アドレスを一つしか持たないようにすることで、2のべき乗の大きさである32バイトに丁度収まるようになっていく。8分木ノードにし、経路アドレスを1ノードで一つだけ保持するようにすれば、64バイトの大きさに収まる上に、4バイトの余裕ができ、この領域は他の情報を入れるのに使用できる。さらに1つにまとめる2分木の数を増やせば(すなわち、2のp乗分木のpの値を大きくとれば)、2のべき乗の大きさに対し、1つのノードの大きさをさらに小さくできる。

【0119】このように、1ノードの大きさを2のべき乗の大きさに収めることで、ハードウェアの構成を非常に簡単にできる。ハードウェアの構成を簡単にできる例を以下に示す。

【0120】1つ目の例として、4分木ノードが32バイトに収まる場合、メモリを複数バンクで構成していた場合でも1ノードのメモリ領域がバンク境界にまたがることはないこと、メモリにダイナミックRAMを使用していた場合でも1ノードのメモリ領域がRowアドレス境界にまたがることはないこと。

【0121】2つ目の例として、ノード内の各要素のメモリ上のアドレス(以下、メモリアドレスと呼ぶ)を求めるときに、そのノードへのポインタとそのポインタからのオフセットの足し算でなく、メモリアドレスの上位ビットはポインタの値にし、下位ビットをオフセットにすればよいこと。例えば4分木ノードが32バイトに収まる場合、あるノード内のある要素のメモリアドレスを生

成するには、そのノードへのポインタをメモリアドレスの2の5乗ビット以上に割り付け、ノード内のその要素へのオフセットをメモリアドレスの2の0乗ビットから2の4乗ビットに割り付ければよい。

【0122】3つ目の例として、例えば4分木ノードが32バイトに収まる場合、各ノードに保持する次のノードへのポインタとして、(次のノードの先頭のバイトアドレス) $\div 32$ 、という値を保持すればよく、1ノード内で1ポインタあたり5ビットずつデータ量を減らせることが挙げられる。

【0123】次に、例えば、図1の経路検索処理回路200を専用LSIで構成する場合のノードの読み込み方式について説明する。4、8、16、...分木ノードにした場合に1つのノードが大きくなり、検索処理時にテーブルを保持するメモリから専用LSI内に1つのノードを全て読み込むと、ノードを大きくするに従い読み込み時間が伸び、性能低下要因となる、という問題があるが、この問題は、ノードを大きくしたときに、1つのノード全てを読み込まず一部だけを読み込む、という方で回避する。この方法について、図20を用いて説明する。

【0124】図20に4分木の場合の例を示す。既に図19を用いて説明したように、経路アドレスマスク長 m ビットの4分木ノードは、図12(b)で示した、経路アドレスの第 m ビット目の値が0の場合に対応する2分木ノードと、1の場合に対応する2分木ノードを併せた構造になっている。従って、受信パケットの経路アドレスの第 m ビット目の値にしたがって、対応する方の2分木ノードの部分だけを読み込むことにより、ノードの大きさが大きくなっても2分木ノードの場合と同じデータ量を読み込むようにする。このとき、図19で示した、1つにまとめられる2分木ノード間で一つだけ保持する要素であるノードの経路アドレス(ワード124、125)は、受信パケットの経路アドレスの m ビット目の値に係わらず読み込むようにする。

【0125】さらに、受信パケットの経路アドレスの第 $(m+1)$ ビット目の値に従い、2分木ノードで2つ存在した次ノードへのポインタの内、一方だけを読み込むようにすることにより、読み込むデータ量をさらに少なくする。

【0126】この方法は2分木検索方式にも使用できる。例えば第 m ビット目の2分木ノードの場合には、受信パケットの経路アドレスの第 m ビット目の値に従い、2つの次ノードへのポインタの内、一方だけを読み込むようにする。

【0127】上記方法を全て行い、結局、このノードの経路アドレスマスク長を m とした場合、宛先アドレスの第 m 、 $m+1$ ビットの値が00か、01か、10か、11かに従い、それぞれ(ワード120→ワード124→ワード125→ワード126)、(ワード121→ワード

124→ワード125→ワード126)、(ワード122→ワード124→ワード125→ワード127)、(ワード123→ワード124→ワード125→ワード127)の順にデータを読み込むようにする(図20(b))。

【0128】このように、あるノードの一部分だけを読み込むためには、このノードの経路アドレスマスク長 m を知る必要があり、このノードの経路アドレスマスク長 m は、1ノードのデータ読み込みの最初に読み込むか、この情報を1つ前のノードに移し、一つ前のノードのデータ読み込み時に読み込む必要がある。ノードのマスク長 m を1ノードのデータ読み込みの最初に読み込む方法は、例えば、図1の経路検索処理回路200を専用LSIで構成する場合、宛先の第 m ビット目の値の抽出のための専用LSI内のゲートディレイ、及び、次に読み込む部分のメモリアドレスをメモリに出力してからメモリからのデータを専用LSI内に読み込むまでの時間であるメモリリードレイテンシだけ、マスク長 m を読み込んでから次に読み込む部分を選択して読み込むまで時間が空いてしまうので、ノードの一部だけを読むことによる性能向上効果が少ない。ノードのマスク長 m を1つ前のノードに移し、一つ前のノードのデータの読み込み時に読み込む方が性能向上効果がある。

【0129】さらに、ノードの経路アドレスマスク長 m を1つ前のノードに移す場合、1ノードのデータを読み込む順序を、1番目に次のノードの経路アドレスマスク長 m 、次のノードへのポインタ、次にノードの経路アドレス、送信ポート情報へのポインタの順にすることにより、次のノードの最初に読み込む部分のメモリアドレスが最も早く計算できるようにする。

【0130】次のノードへのポインタは、次のノードのメモリ領域の先頭部分を指しており、次のノードの先頭から最初に読み込む部分までのメモリアドレスのオフセットは、次のノードのマスク長 m を読み込み、受信パケットの経路アドレスの該当ビット位置の値を検査することにより、得られる。

【0131】次に、1ノード内で、条件によっては、読み込む必要が無い要素を、その条件に従って読み込まないようにすることで、読み込みの時間を削減し、高速化を図る方法について図21を用いて説明する。

【0132】図21に4分木の場合の例を示す。2分木構造の説明で述べたように一般に木構造においては、全ノードにエントリが割り付けられているわけではなく、枝の分岐の個所では、エントリが割り付けられていなくてもノードを設ける必要がある。図21に示すように、ノードデータの最初に読み込むワード120、121、122、123内のフラグ内のエントリ有りフラグの値に従い、エントリが割り付けられていないノードでは、送信ポート情報へのポインタを読み込まないようにすることで、読み込み時間の短縮を図ることができる。エン

トリ有りフラグは1ビットで表現できるので、この情報を読み込むことによる読み込み時間の増大は小さい。

【0133】この方法では、このノードのマスク長を m とすると、宛先アドレスの第 m 、 $m+1$ ビットの値が00で、ワード120を読み、フラグ00から、4分木を構成する0番目の2分木に経路情報が無いと判った場合、ワード124、125だけを読めば良く、経路情報が有ると判った場合にだけ、図21(b)に示すようにワード124→ワード125→ワード126の順に読めば良い。宛先アドレスの第 m 、 $m+1$ ビットの値が01、10、11の場合も同様である。

【0134】以上、2分木検索方式(方式1)、2の p 乗分木検索方式(方式2)、およびマスク長 m ビットのノードをメモリ上に展開する方法(方式3)という各検索方式について説明した。次に、図1の経路検索処理回路200が上記の方式を用いて経路検索処理を行う際のフローチャートについて図22を用いて説明する。図22のフローチャートでは、方式2と方式3とを組み合わせた場合の例を示す。この例では、経路アドレスの第0ビットから第 $(m-1)$ ビットの値に従って、メモリ上の決まった位置に展開された2の p 乗分木ノードの一つを選択する。以下ではこのように選択し、検索の最初に読み込む2の p 乗分木ノードを初段ノードと呼ぶ。第 m ビット以降は経路アドレスを p ビットずつ検索し、2の p 乗分木を検索する。図1のメモリ300には、上記の方式2および方式3の検索方式に従った2の p 乗分木ノードデータ、および送信ポート情報が格納されているとする。

【0135】なお、以下では方式2と方式3とを組み合わせた例について説明するが、方式1と方式3とを組み合わせた場合でも同様なフローチャートに従うことで実現可能である。また、このフローチャートに従うことにより、ソフトウェアでもハードウェアでも経路検索処理を実現することができる。ソフトウェアで実現する場合、図1の経路検索処理回路200にはCPUを用いればよい。またハードウェアで実現する場合、図1の経路検索処理回路200を専用LSIで構成すればよい。

【0136】図22の処理810は木構造検索処理であり、また、処理811は送信ポート情報処理である。まず、木構造検索処理810について説明する。

【0137】図1の経路検索処理回路200は、受信パケットの宛先マルチキャスト・グループIDおよび送信元IPアドレスを受け取ると、これらの値から図5を用いて説明した経路アドレスを生成し、この経路アドレスと初段ノードのマスク長 m の値から初段ノードへのポインタを生成し、このポインタと、経路アドレスの第 m ビットから第 $(m+p-1)$ ビットの値(以下、検査ビット値と呼ぶ)に従ってメモリ300に格納されている初段ノードの読み込みアドレスを生成し、メモリ300から該初段ノードの一部を読み込む(図22の800)。

【0138】次に、図1の経路検索処理回路200は、

受信パケットの経路アドレスにノードのマスク長だけ上位ビットから有効とするマスクをかけたものと、ノードの経路アドレスと比較し(図22の801)、不一致の場合は木構造検索処理810を終了する(図22の809)。一致する場合は図22の処理802に進む。

【0139】次に、図22の処理802、および803について説明する。これらの処理は、最長一致検索を実現するための経路情報の更新処理である。経路情報としては、図1のメモリ300が出力するノードデータの内、エントリ有りフラグと、送信ポート情報へのポインタがある。図1の経路検索処理回路200は、ノードデータ内のエントリ有りフラグを検査し、エントリ有りフラグの値が1の場合(図22の812)のみ、読み込んだノード内の新たな経路情報をレジスタに保持する(図22の803)。エントリ有りフラグの値が0の場合は更新処理を行わない(図22の813)。

【0140】次に、図1の経路検索処理回路200は、図1のメモリ300が出力するノードデータの中の次ノードへのポインタがNULLかどうかを判定し(図22の804)、NULLの場合は木構造検索処理810を終了する。NULLでない場合はそのポインタと、新たな検査ビットの値に従ってメモリ300に格納されている初段ノードの読み込みアドレスを生成し、メモリ300から該ノードデータを読み込む(図22の805)。

【0141】以下、上記の処理を繰り返すことにより、2の p 乗分木方式の経路検索を行うことができる。

【0142】次に、図22の送信ポート情報処理811について説明する。木構造検索処理の結果、図1の経路検索処理回路200内のレジスタには、経路情報として、上記のエントリ有りフラグと送信ポート情報へのポインタが保持されている。図1の経路検索処理回路200は、まず、レジスタに保持されているエントリ有りフラグを調べ(図22の806)、その値が0の場合は経路検索処理を終了し、転送処理部504へ検索結果無しという通知をする。エントリ有りフラグの値が1の場合は、検索の結果、あるエントリに一致したことになるため、送信ポート情報へのポインタを用いてメモリ300から送信ポート情報を読み出す(図22の807)。この送信ポート情報からパケットを送信すべき1つまたは複数のポート番号を生成して図1の転送処理部504へ送信し(図22の808)、すべてのポート番号を送信し終わると経路検索処理を終了する。

【0143】次に、本発明の一実施例として、図22で説明した検索方式をハードウェアで実現する場合の構成例について、図23、図24を用いて説明する。

【0144】図23に経路検索処理部400をハードウェアで構成した場合の構成例を示す。経路検索処理回路200は、木構造検索回路201と、読み込みアドレス生成回路202と、メモリ制御回路206と、送信ポート情報処理回路203と、経路検索処理制御回路204

とからなる。

【0145】木構造検索回路201は、メモリ300に格納された2のp乗分木構造を検索し、次に読み込むべきノードのポインタの生成、受信パケットの経路アドレスの検査ビット値の抽出、木構造検索の終了判定、検索結果である経路情報の候補の更新を行う。また、読み込みアドレス生成回路202は、木構造検索処理時には、木構造検索回路201から出力される読み込むべきノードへのポインタ、および検査ビット値、および経路検索処理制御回路204から出力されるタイミング信号に従い、図20を用いて説明した順序でノード内の読み込みワードのメモリアドレスを生成してメモリ制御回路206に送信する。また、送信ポート情報処理時には、木構造検索回路201内のレジスタに検索結果として保持されている送信ポート情報へのポインタから送信ポート情報のメモリアドレスを生成する。また、メモリ制御回路206は、上記のメモリアドレスと、経路検索処理制御回路204から出力されるタイミング信号に従い、メモリ制御信号を生成する。また、送信ポート情報処理回路203は、メモリ300から読み出された送信ポート情報からパケットを送信すべき1つまたは複数のポート番号を生成して転送処理部504へ送信する。また、経路検索処理制御回路204は、経路検索処理回路200全体の制御(各回路の動作タイミングおよび動作状態管理など)を行う。

【0146】上記の経路検索処理回路200を専用LSIで実現する場合、メモリ300は上記専用LSIの外部に備えても良いし、専用LSI内の内蔵メモリを使用しても良い。メモリ300として専用LSIの内蔵メモリ使用すると、外部メモリを使用する場合に比べてメモリアクセス時間が短くなり、メモリ300内に格納されたルーティングテーブルの検索処理を高速に行うことができる。

【0147】次に、経路検索処理部400の動作について図23を用いて説明する。また、木構造検索回路201の詳細動作については、図24を用いて後述する。

【0148】木構造検索回路201は、転送処理部504から受信パケットの宛先マルチキャスト・グループIDおよび送信元IPアドレスを受け取ると、これらの値から経路アドレスを生成し、この経路アドレスとノードのマスク長の値から次ノードへのポインタを生成して、読み込みアドレス生成回路202に転送する。また、木構造検索回路201は、ノードのマスク長で示される経路アドレスの検査ビット位置の値(検査ビット値)を抽出して、読み込みアドレス生成回路202に転送する。

【0149】読み込みアドレス生成回路202はこのノードへのポインタと、検査ビット値と、経路検索処理制御回路204からのタイミング信号を用いて、図20を用いて説明した順序でノード内の読み込みワードのメモリアドレスを生成し、メモリ制御回路206に送信す

る。メモリ制御回路206は上記メモリアドレスと経路検索処理制御回路からのタイミング信号を用いてメモリ制御信号を生成し、メモリ300へ出力する。上記のメモリ制御信号を受信したメモリ300は、対応するノードデータを信号線214を用いて木構造検索回路へ転送する。

【0150】木構造検索回路201はこのノードデータを用いて、図22の処理801、802、803、804、805を行う。これらの処理の詳細は図24で後述する。図22の処理801および804に対応する判定処理において木構造検索を終了すると判定した場合は、木構造検索終了信号を経路検索処理制御回路204へ出力し、経路検索処理制御回路204は送信ポート情報読み込み処理を開始する。終了と判定されない場合は、終了と判定されるまで図22の処理801、802、803、804、805を繰り返す。

【0151】次に、経路検索処理制御回路204は、木構造検索回路201内に保持された経路情報(エントリ有りフラグ、送信ポート情報へのポインタ)のうちのエントリ有りフラグを調べ、その値が0の場合は経路検索処理を終了し、転送処理部504へ検索結果無しという通知をする。エントリ有りフラグの値が1の場合は、送信ポート情報読み込みおよび送信ポート生成処理制御を開始し、木構造検索回路201を制御して、送信ポート情報へのポインタを読み込みアドレス生成回路202へ出力させる。読み込みアドレス生成回路202はこの送信ポート情報へのポインタと、経路検索処理回路204からのタイミング信号を用いて、読み出すべき送信ポート情報が格納されているメモリアドレスを生成し、メモリ制御回路206へ送信し、メモリ制御回路206は、上記のメモリアドレスと、経路検索処理制御回路204から出力されるタイミング信号に従い、メモリ制御信号を生成してメモリ300へ出力し、この制御信号を受信したメモリ300から信号線214上に送信ポート情報が出力される。送信ポート情報処理回路203はこの送信ポート情報を取り込み、パケットを送信すべき1つまたは複数のポート番号を生成して転送処理部504へ送信する。

【0152】送信ポート情報の構造については、複数の送信ポート番号をポインタでつないだリスト構造や、送信ポート番号のビットマップ化などが考えられるが、リスト構造の場合、送信ポート番号が多いとメモリアクセス回数が増加し、結果として検索処理が遅くなる問題がある。送信ポート番号をビットマップ化する方法は、そのビットマップをメモリ300からバーストリードして送信ポート情報処理回路203内に保持しておき、そのビットマップをデコードして送信ポート番号を生成することにより、メモリアクセス回数を減らすことができる。

【0153】すべての送信ポート番号を出力した後で、

経路検索処理制御回路204は経路検索処理を終了し、次のバケット処理の制御を行う。

【0154】次に、図23の木構造検索回路201の詳細を図24を用いて説明する。

【0155】まず、図22の初段ノードリード処理800および次ノードリード処理805に対応する処理について、図24を用いて説明する。転送処理部504より信号線205を用いて送信される宛先マルチキャスト・グループアドレスおよび送信元IPアドレスは、経路アドレス生成回路207内に保持される。経路アドレス生成回路207はこれらの情報から、図5(b)で説明した60ビットの経路アドレス33を生成し、初段ノードへのポインタ生成回路208、検査ビット抽出回路209、マスク処理回路210へ出力する。

【0156】初段ノードへのポインタ生成回路208は、予め初段ノードマスク長レジスタ229に設定されている値 m に従い、受信バケットの経路アドレスの第0ビットから第 $(m-1)$ の上位 m ビットの値を抽出し、この m ビットの値に従って初段ノードへのポインタを生成してポインタセクタ211へ出力する。ポインタセクタ211は、信号線226によって図23の経路検索処理制御回路204から送信される初段ノード読み込み／初段ノード以外のノード読み込み選択信号に従い、初段ノードの読み込み時は、上記の初段ノードへのポインタを選択して図23の読み込みアドレス生成回路202へ出力する。また、ポインタセクタ211は、初段ノード以外のノードの読み込み時は、次ノードへのポインタレジスタ215に保持されている次ノードへのポインタを選択して読み込みアドレス生成回路202へ出力する。

【0157】また、上記の処理と並行して、検査ビット抽出回路209は、初段ノード読み込み時には、検査ビット位置セクタ212において選択されて出力される初段ノードマスク長レジスタ229の設定値 m に従い、経路アドレスの第 m ビットから第 $(m+p-1)$ ビットまでの p ビットの検査ビット値を抽出して図23の読み込みアドレス生成回路202へ出力する。また、初段ノード以外のノードの読み込み時は、後述する次ノードマスク長レジスタ216に保持されている次ノードのマスク長 $m1$ が、検査ビット位置セクタ212において選択されて検査ビット抽出回路209に出力され、検査ビット抽出回路209は、この次ノードのマスク長 $m1$ に従い、経路アドレスの第 $m1$ ビットから第 $(m1+p-1)$ ビットまでの p ビットの検査ビット値を抽出して図23の読み込みアドレス生成回路202へ出力する。

【0158】図23の読み込みアドレス生成回路202およびメモリ制御回路206は、上記のノードへのポインタと p ビットの検査ビット値を用いて、図20を用いて説明した順序でノード内の読み込みワードのメモリアドレスおよびメモリ制御信号を生成し、メモリ300は

入力されたメモリアドレスに格納されているノード内の読み込みワードを信号線214に出力する。

【0159】信号線214のビット幅が32ビットの場合、図20で説明した方式を採用すると一回のノード(の一部)の読み込みワード数は4ワードとなり、これらのデータは図20の表に示した順番で信号線214に出力され、各レジスタ(次ノードへのポインタレジスタ215、次ノードマスク長レジスタ216、経路アドレス上位レジスタ217、経路アドレス下位レジスタ222、フラグレジスタ224、送信ポート情報ポインタレジスタ230)に保持される。各レジスタの保持タイミングは経路検索処理制御回路204からの制御信号(図示していない)により制御される。次ノードへのポインタレジスタ215には次ノードへのポインタが保持され、次ノードマスク長レジスタ216には次ノードのマスク長が保持され、経路アドレス上位レジスタ217にはノードの経路アドレスの上位28ビットが、経路アドレス下位レジスタ222には経路アドレスの下位32ビットがそれぞれ保持され、フラグレジスタ224にはエントリ有りフラグが保持され、送信ポート情報ポインタレジスタ230には送信ポート情報へのポインタが保持される。

【0160】次に、図22の経路アドレス一致比較処理801に対応する処理について、図24を用いて説明する。初段ノード読み込み時には、マスク長セクタ228を介して初段ノードマスク長レジスタ229の設定値がマスク処理回路210に入力する。

【0161】初段ノード以外のノード読み込み時には、マスク長セクタ228を介してマスク長レジスタ227に保持されているノードのマスク長がマスク処理回路210に入力する。このマスク長レジスタ227の値は、一つ前に読み込んだノードに格納されている次ノードのマスク長の値であり、次ノードマスク長レジスタ216に保持されていたものである。次ノードマスク長レジスタ216の値は、ノードの読み込み毎に更新されるため、更新される前に、現在読み込んでいるノード(以下、現ノードと呼ぶ)の経路アドレス一致比較処理に使用する現ノードのマスク長をマスク長レジスタ227に保持しておく。

【0162】マスク処理回路210は、これらのマスク長だけ上位ビットから有効とするマスクを生成し、このマスクと、経路アドレス生成回路207から出力される受信バケットの経路アドレスとの論理積をとり、その結果(以下、マスクした受信バケットの経路アドレスと呼ぶ)を一致比較回路213に出力する。また、経路アドレス上位レジスタ217、経路アドレス下位レジスタ222に保持されているノードの経路アドレスの上位28ビット、下位32ビットはこの順に連結され、ノードの経路アドレスとして一致比較回路213に入力される。一致比較回路213は、このノードの経路アドレスと、

上記で説明した、マスクされた受信バケットの経路アドレスとを比較し(以下経路アドレス比較と呼ぶ)、その結果が不一致となる場合には不一致信号を信号線218を用いて木検索終了判定回路219に転送する。木検索終了判定回路219は上記不一致信号を受信し、木構造検索終了信号を図23の経路検索処理制御回路204に送信する。

【0163】また、マスク長レジスタ227の値は、上記の経路アドレス比較が行われた後、次ノードマスク長レジスタ216に保持されている次ノードのマスク長の値によって更新され、次のノードの読み込み時の経路アドレス比較に使用される。

【0164】次に、図22の処理802、および803に対応する処理について、図24を用いて説明する。

【0165】一致比較回路213における経路アドレス比較の結果一致した場合、一致比較回路213は信号線220を用いて一致信号を更新判定回路221に出力する。この一致信号が更新判定回路221に入力され、かつ、フラグレジスタ224に保持されているエントリ有りフラグの値が1の場合(図25の812)のみ、更新判定回路221は更新信号をフラグレジスタ223および送信ポート情報へのポインタレジスタ225に出力する。更新信号を受信したフラグレジスタ223は、フラグレジスタ224に保持されているエントリ有りフラグを新たに保持し、同じく更新信号を受信した送信ポート情報へのポインタレジスタ225は、送信ポート情報ポインタレジスタ230に保持されている送信ポート情報へのポインタを新たに保持する(図22の経路情報更新処理803)。エントリ有りフラグの値が0の場合は更新判定回路221は更新信号を送信しないので、フラグレジスタ223および送信ポート情報へのポインタレジスタ225は更新処理を行わない(図22の813)。

【0166】経路情報として上記のフラグおよび、送信ポート情報へのポインタ以外の情報を追加する必要がある場合は、木構造のノード内にそれらの情報を追加し、それらの情報を保持、更新するレジスタを新たに追加すればよい。

【0167】次に、図22の次ノードへのポインタがNULLかどうかの判定処理804に対応する処理について、図24を用いて説明する。次ノードへのポインタレジスタ215内に保持されている次ノードへのポインタは木検索終了判定回路219に入力される。この次ノードへのポインタがNULLの場合、木検索終了判定回路219は図23の経路検索処理制御回路204に木検索終了信号を出力する。

【0168】以上、図23の木構造検索処理回路201の詳細動作を図24を用いて説明したが、ハードウェアで構成するため、図22の木構造検索処理の中の801、802、803、804、805の各処理は逐次処理をする必要はなく、各処理に必要なデータが各レジス

タ215、216、217、222、224、230に保持された後に、各処理を開始すればよく、上記の各処理を並列処理を行うことにより、高速に木構造の検索を行うことができる。

【0169】

【発明の効果】マルチキャスト経路検索において検索のキーとなる宛先マルチキャスト・グループアドレスと送信元IPアドレスをこの順に連結して一つの経路アドレスとして定義し、この経路アドレスのビットパターンに従い、マルチキャストルーティングテーブルを2分木に構成することにより、エントリ数が増大しても検索時間は増加しない検索処理を行うことができる。

【0170】また、検索木を構成する各ノードを、従来の方法である2分木ノードから4分木、8分木、あるいはそれ以上と、枝別れの数を2のべき乗で増やすことにより、1つのノードで1ビットでなく、連続する2ビット、3ビット、あるいはそれ以上のビット数を同時に検査でき、検索終了までに辿るノードの数が減り、経路検索処理の高速化を図る効果がある。

【0171】図25に、エントリ数に対する従来技術2のハッシュ方式のマルチキャスト経路検索の処理性能900、および本発明の2分木方式のマルチキャスト経路検索の処理性能901、2のp乗分木方式のマルチキャスト経路検索の処理性能(例として、4分木方式の性能902、8分木方式の性能903)を示す。ハッシュ方式では、エントリ数に反比例して処理性能が劣化するのに対し、2分木方式および2のp乗分木方式ではエントリ数に依存せず一定となる。また、2分木方式に対し、4分木方式は2倍の性能を実現でき、また、8分木方式では3倍の性能を実現できる。

【0172】さらに、上記の経路アドレスとして、マルチキャスト・グループアドレス自体ではなく、マルチキャスト・グループアドレスの下位28ビットのマルチキャスト・グループIDと送信元IPアドレスをこの順に連結したものを採用し、この経路アドレスのビットパターンに従い、マルチキャストルーティングテーブルを2分木あるいは2のp乗分木に構成することにより、検索木の深さを決定する経路アドレスのビット数を4ビット減らし、ノードを渡る最大回数を減らすことができる。

【0173】また、4分木、8分木、あるいは一般に2のp乗分木を構成するときに、1つの2分木ノードと、その直下につながる(p-1)段分の合計(2のp乗-1)個分の2分木ノードを一つの2のp乗分木ノードにまとめ、まとめられる最下段の2の(p-1)乗個の2分木ノードに、それより上段のノードに割り付けられていた経路データを埋め込むことにより、まとめる前には2分木ノード換算で(2のp乗-1)個分のメモリ量だったものを、2の(p-1)乗個分のメモリ量に減らす効果がある。

【0174】さらに、2分木ノードを2のp乗分木にま

とめる際、送信ポート情報を検索木のノードから分離して記憶手段の別領域に格納し、検索木のノードには上記送信ポート情報の格納領域へのポインタを保持することで、エントリが割り当てられている2分木ノードと、エントリが割り当てられていない2分木ノードのデータ量を同じにし、これらの2分木ノードをまとめて2の p 乗分木ノードを構成する際に、2の p 乗分木ノードのメモリ量の削減と、各2の p 乗分木ノードのメモリ量の統一を図ることができる

また、この2分木を複数個併せた形で構成した4分木、8分木、あるいはそれ以上の枝別れ数のノードを、検索のために読むときにノード全てを読むのではなく、必要な部分のみを読むようにすることにより、ノードが大きくなることによるデータの読み込み時間の増大を防ぎ、経路検索処理の高速化を図る効果がある。

【0175】また、各ノードに、そのノード自体のマスク長ではなく、そのノードのすぐ下につながるノードのマスク長を格納することにより、ノードのマスク長をノードのデータを読む前に知ることができ、ノードのデータを読む前に、ノードのどの部分を読み込めば良いかが分かり、必要な部分のみを読むことができるようになり、ノードが大きくなることによるデータの読み込み時間の増大を防ぎ、経路検索処理の高速化を図る効果がある。

【0176】また、各ノードの先頭に、そのノードにエントリが割り付けられているか否かを示すフラグを設け、最初に、このフラグを読み込み、エントリが割り付けられていないノードでは、経路情報を読み込まないようにすることにより、データの読み込み時間が短縮され、経路検索処理の高速化を図る効果がある。

【0177】また、経路アドレスの上位数ビット分、ノードをメモリ上の決まった位置に展開し、受信パケットの経路アドレスの上位数ビットに従い、ノードが格納されている位置を直接にアクセスすることは、検索処理時間を無くし、経路検索処理の高速化を図る効果がある。

【0178】また、上記の高速化手段をハードウェアで構成することにより、経路検索処理の高速化を図る効果がある。

【図面の簡単な説明】

【図1】本発明の一実施例であるパケット中継装置の例であるルータのブロック図。

【図2】一般的なネットワークの構成例を示す図。

【図3】IPアドレスの構成、およびクラスDのIPアドレスの構成、およびIPパケットのパケットヘッダフォーマット。

【図4】一般的なルータの構成例を示す図。

【図5】本発明で定義する経路アドレスの構成図。

【図6】経路アドレス長3ビットの場合の全てのノードがある2分木。

【図7】経路アドレス長3ビットの場合のマルチキャスト

ルーティングテーブルの例。

【図8】図6においてエントリが割り付けられておらず、かつ、エントリ付きのノードへの途中経路にもなっていないノードを取り除いた木。

【図9】経路アドレス32ビットの場合のマルチキャストルーティングテーブルの例。

【図10】図9に示したマルチキャストルーティングテーブルに対応する木。

【図11】図10において枝別れもエントリの割り付けも無いノードを取り除いた木。

【図12】2分木ノードのデータ構造を示す図。

【図13】2分木から8分木への変形時に一つの8分木ノードにまとめられる2分木ノードを囲んだ木。

【図14】マスク長が0ビットから始まらないようにビット位置を区切った場合に最初のノードをメモリ上に展開することにより区切り位置までのビットの検索を行う木。

【図15】先頭のさらに多くのビット数の検索時間を省いた木。

【図16】2分木から4分木への変形時に1つの4分木ノードにまとめられる2つの2分木ノードを2つの2分木ノード分につぶした木。

【図17】1つの4分木ノードにまとめられる3つの2分木ノードの、2つの2分木ノードへのつぶし方。

【図18】1つの8分木ノードにまとめられる7つの2分木ノードの、4つの2分木ノードへのつぶし方。

【図19】4分木ノードのデータ構造を示す図。

【図20】ノードの大きさが大きくなったときにノードデータリード時間の増大を防ぐために1つのノード全てを読み込まずに一部だけを読み込む方法を示す図。

【図21】条件によっては読み込む必要が無い要素を、条件に従い読み込まないようにすることで読み込みの時間を削減することで高速化を図る方法を示す図。

【図22】本発明の一実施例である経路検索処理部の検索処理フローチャート。

【図23】本発明の一実施例である経路検索処理部のブロック図。

【図24】本発明の一実施例である、木構造のマルチキャストルーティングテーブル検索を行う木構造検索処理回路のブロック図。

【図25】ハッシュ方式、2分木方式、2の p 乗分木方式の、エントリ数に対する処理性能を示す図。

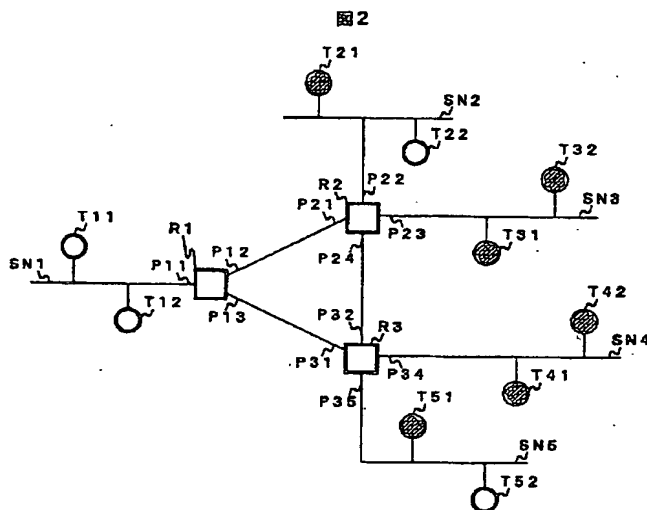
【符号の説明】

T11、T12、T21、T22、T31、T32、T41、T42、T51、T52…ネットワーク上でパケットを送受信する端末、SN1～SN5…各端末が属するサブネットワーク、R1～R3…サブネットワークを相互に接続するルータ、P11、P12、P13、P21、P22、P23、P24、P31、P32、P34、P35…各ルータのポート、1…IPア

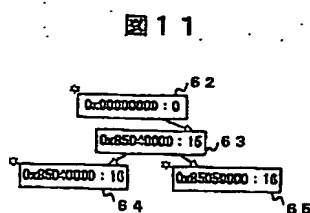
ドレス、2…サブネットワークを識別するサブネットワークアドレス、3…ネットワーク内の端末を識別するホストアドレス、4…クラスDのIPアドレス、5…クラスDであることを示すクラスDのIPアドレスの上位4ビット、6…マルチキャスト・グループID、7…IPパケットヘッダの送信元アドレスフィールド、8…IPパケットヘッダの宛先アドレスフィールド、11-i(=1~N)…ルータの入力ポート、12-i(=1~N)…ルータの出力ポート、13-i(=1~N)…ルーティング処理部、14…スイッチ、15-i(=1~N)…通信制御部、16…管理部、32…マルチキャスト・グループアドレス、送信元IPアドレスの順に連結して生成した経路アドレス、33…マルチキャスト・グループID、送信元IPアドレスの順に連結して生成した経路アドレス、60…経路アドレス長3ビットの場合の受信パケットの経路アドレスの一例、40~54、62~65、70~84、90~95、100~104、130~144…木を構成するノード、61…枝別れもエントリの割り付けも無いノードの長い列、110~113、115…2分木ノードを構成する3

2ビットを単位とするワード、114…送信ポート情報、120~127…4分木ノードを構成する32ビットを単位とするワード、200…経路検索処理回路、201…木構造検索回路、202~204、206…経路検索処理回路を構成する各回路、205…転送処理部、経路検索処理部間の信号線、207~213、215~230…木構造検索回路を構成する各回路および信号線、226…経路検索処理制御回路からの読み込みノード(初段ノードかそれ以外か)選択制御信号線、300…マルチキャストルーティングテーブルが格納されるメモリ、400…経路検索処理回路とメモリから構成される経路検索処理部、500、604、605…ルータに接続されるネットワークの例、501、602、603…通信制御部、502…パケットデータのメモリ格納およびメモリからの読み出しを制御する制御回路、503…パケットデータが格納されるメモリ、504…パケット転送処理部、505…複数のパケット処理部および管理部が接続されるスイッチ、600…パケット処理部、601…ルータ全体の管理を行う管理部、700…ルータ。

【図2】



【図11】

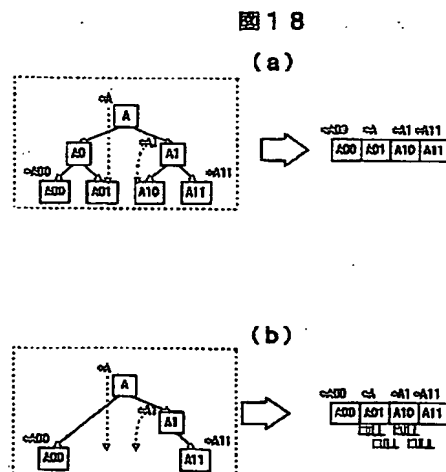


【図7】

図7

経路アドレス	送信元アドレスマスク
000	0
001	3
010	2
100	1
110	3

【図18】



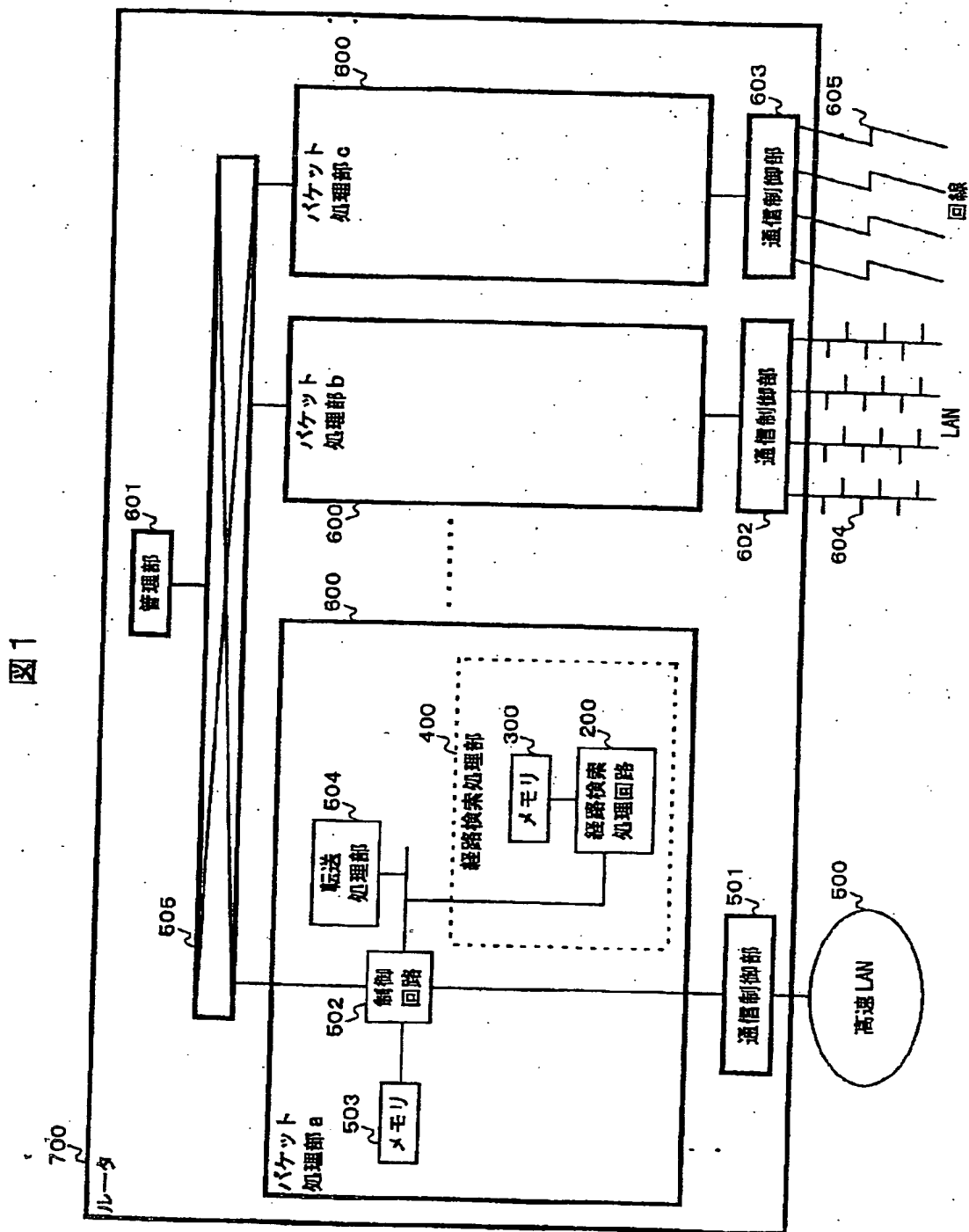
【図9】

図9

マルチキャストルーティングテーブル

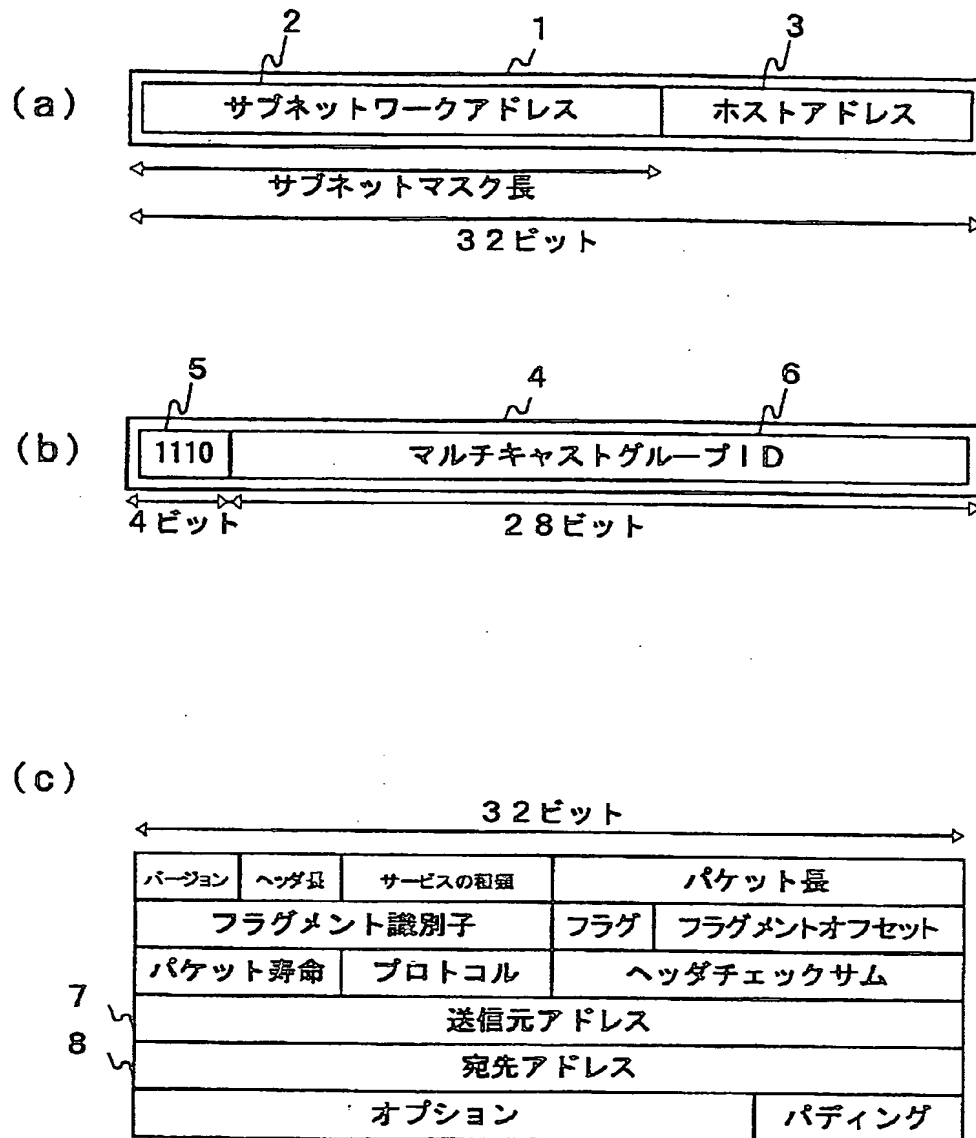
送信元アドレス	送信元アドレスマスク
0x00000000	0
0x85000000	16
0xC0000000	16

【図1】



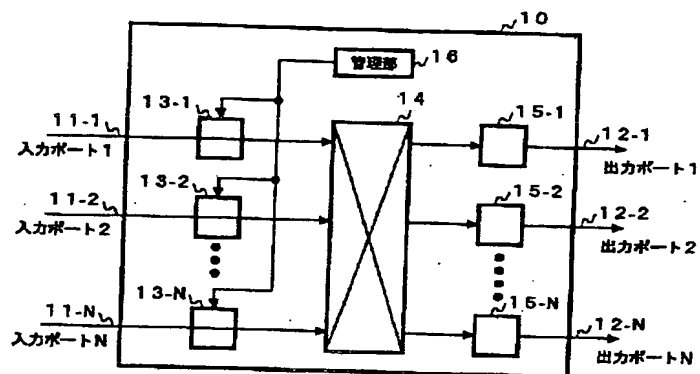
【図3】

図3



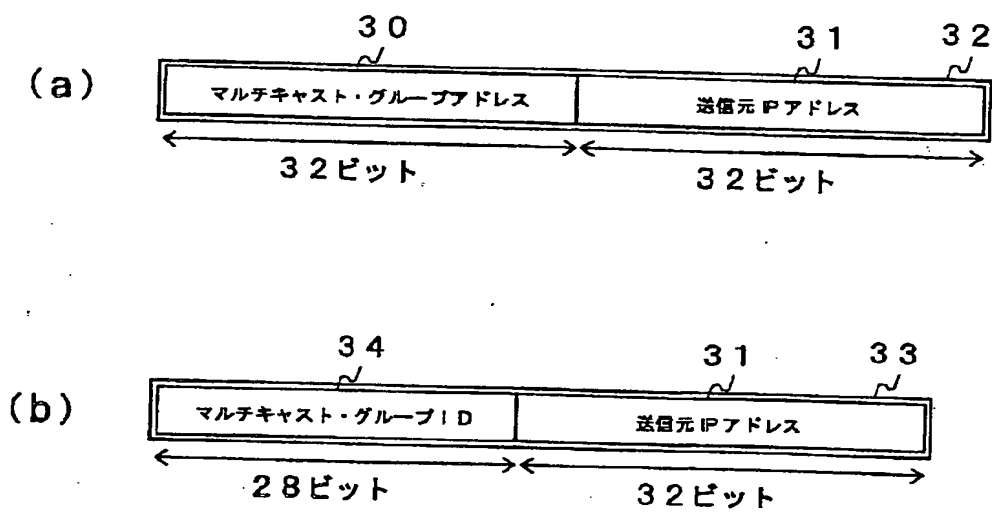
【図4】

図4



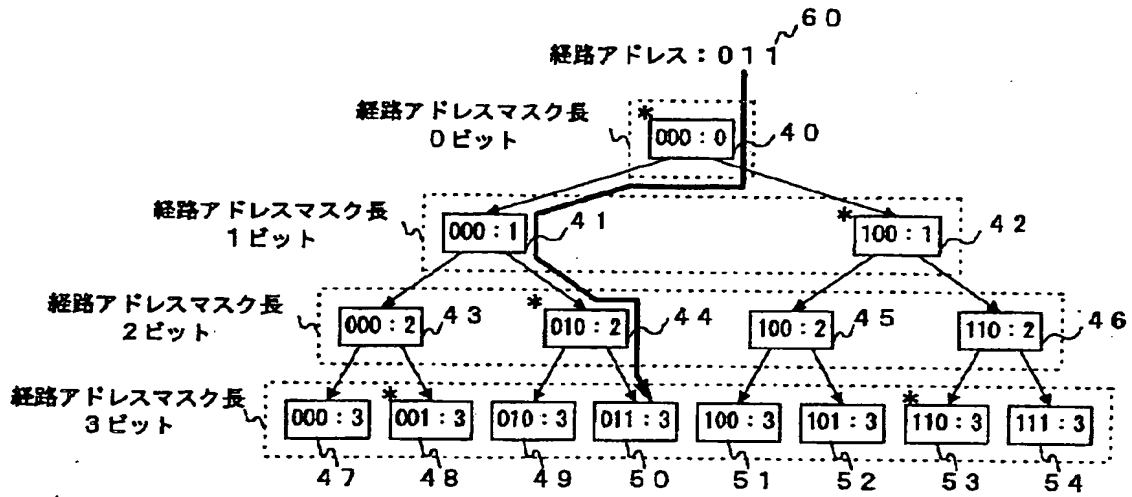
【図5】

図5



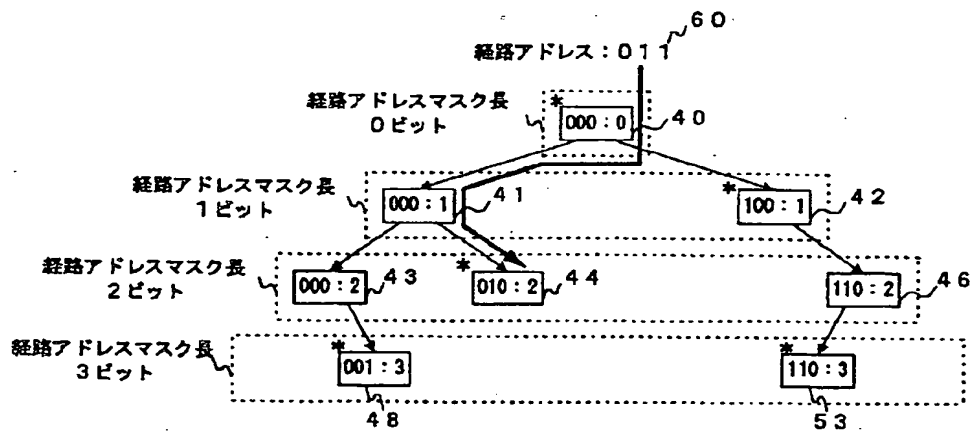
【図6】

図6



【図8】

図8



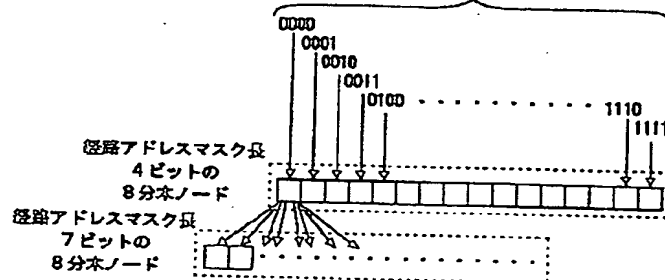
【図 15】

15

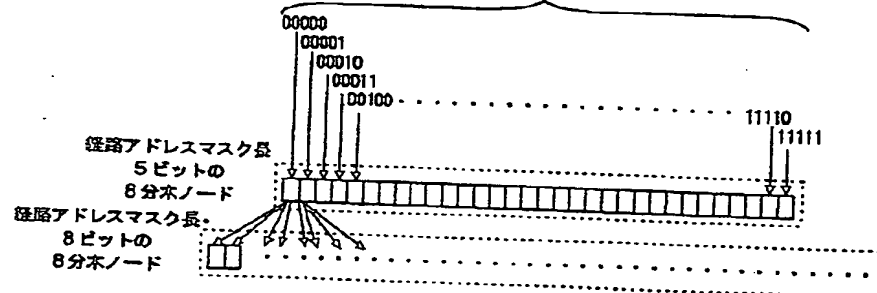


The diagram illustrates two types of address masks. The top part shows a '経路アドレスマスク長 3ビットの 8分本ノード' (3-bit path address mask length, 8-way node). It features a horizontal row of 8 rectangular nodes. Above each node is a 3-bit binary address: 000, 001, 010, 011, 100, 101, 110, and 111. Vertical arrows point from each address down to its corresponding node. The bottom part shows a '経路アドレスマスク長 6ビットの 8分本ノード' (6-bit path address mask length, 8-way node). It features a horizontal row of 2 rectangular nodes. Above each node is a 6-bit binary address: 000000 and 000001. Vertical arrows point from each address down to its corresponding node. Dashed lines indicate the boundaries of the mask lengths.

受信バケットの経路アドレス
第0～3ビットの値により選択

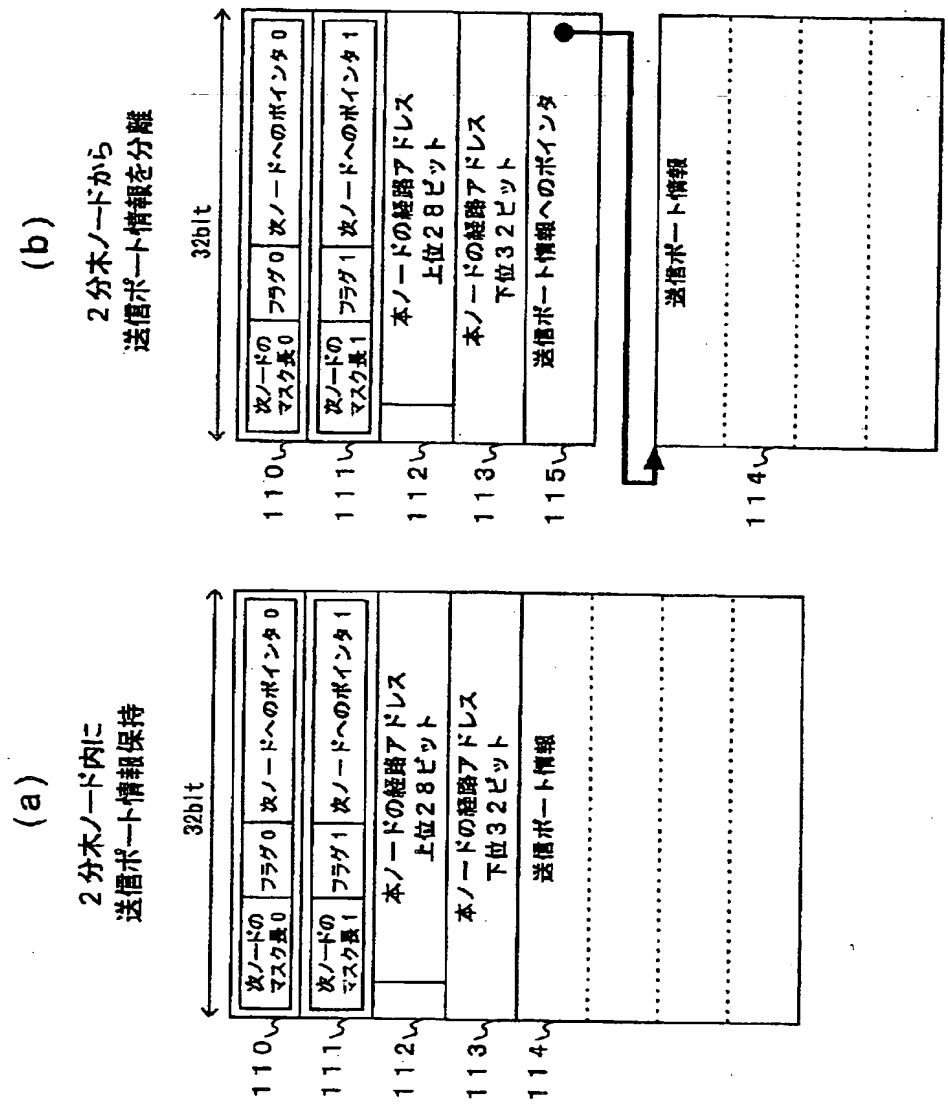


受信バケットの経路アドレス
第0～4ビットの値により選択

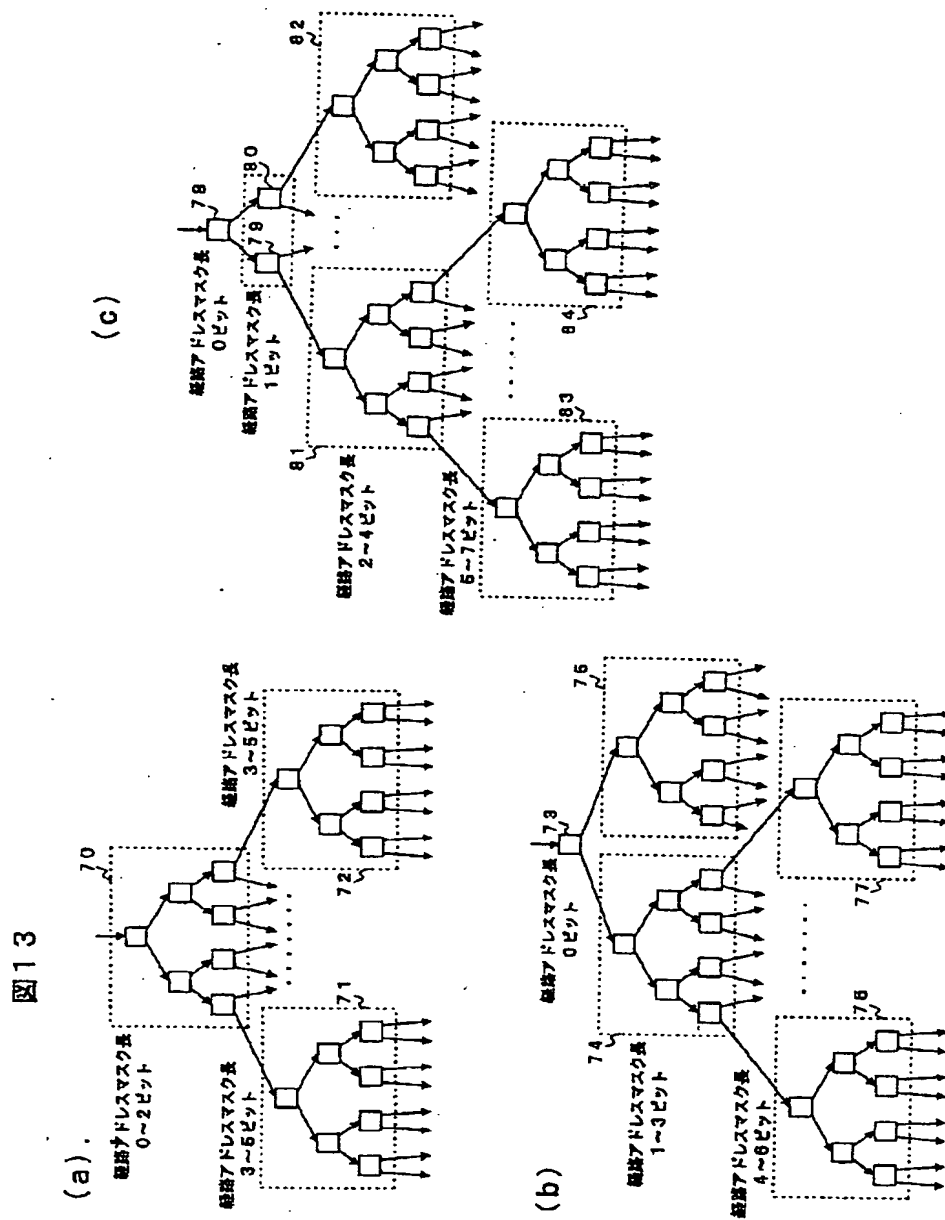


【図12】

図12



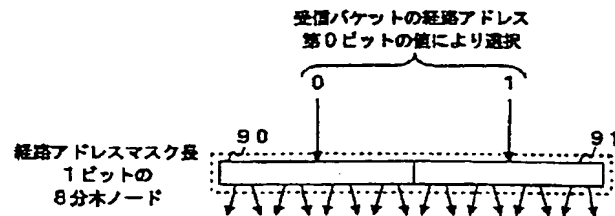
【図13】



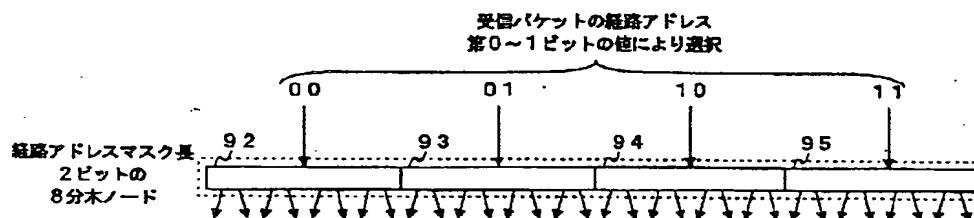
【図14】

図14

(a)



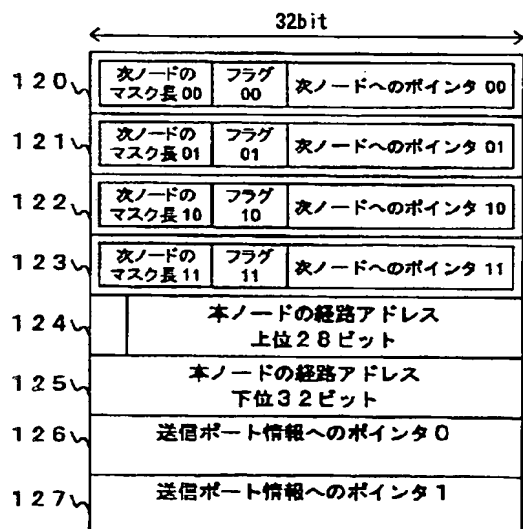
(b)



【図20】

図20

(a)



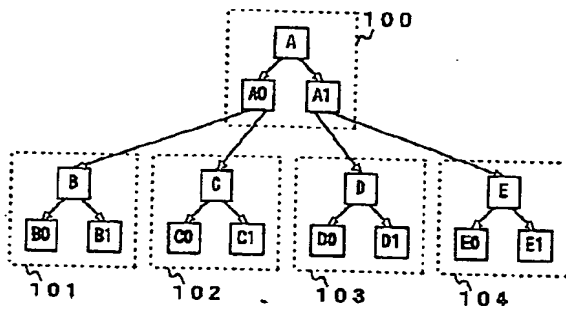
(b)

検査ビットの値	00	01	10	11
読み込ワード	120			
	↓	121		
	↓	↓	122	
	↓	↓	↓	123
	124	124	124	124
	↓	↓	↓	↓
	125	125	125	125
	↓	↓	↓	↓
	126	126		
			127	127

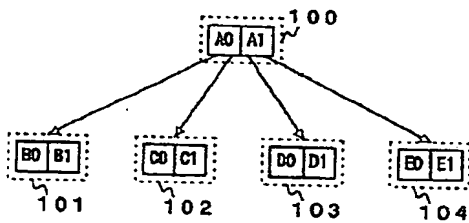
【図16】

図16

(a)

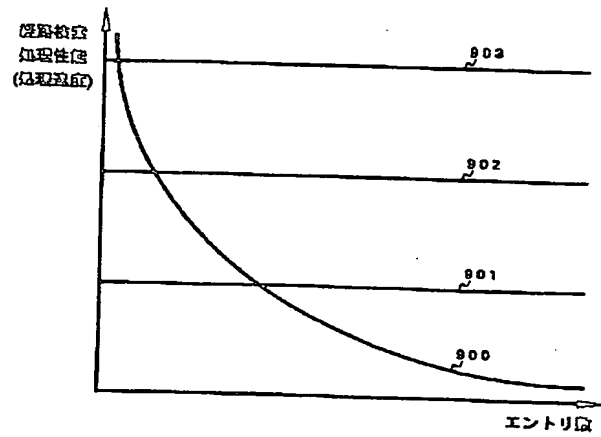


(b)



【図25】

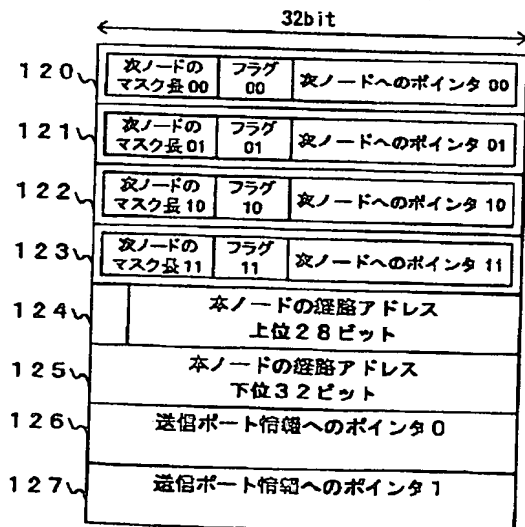
図25



【図21】

図21

(a)

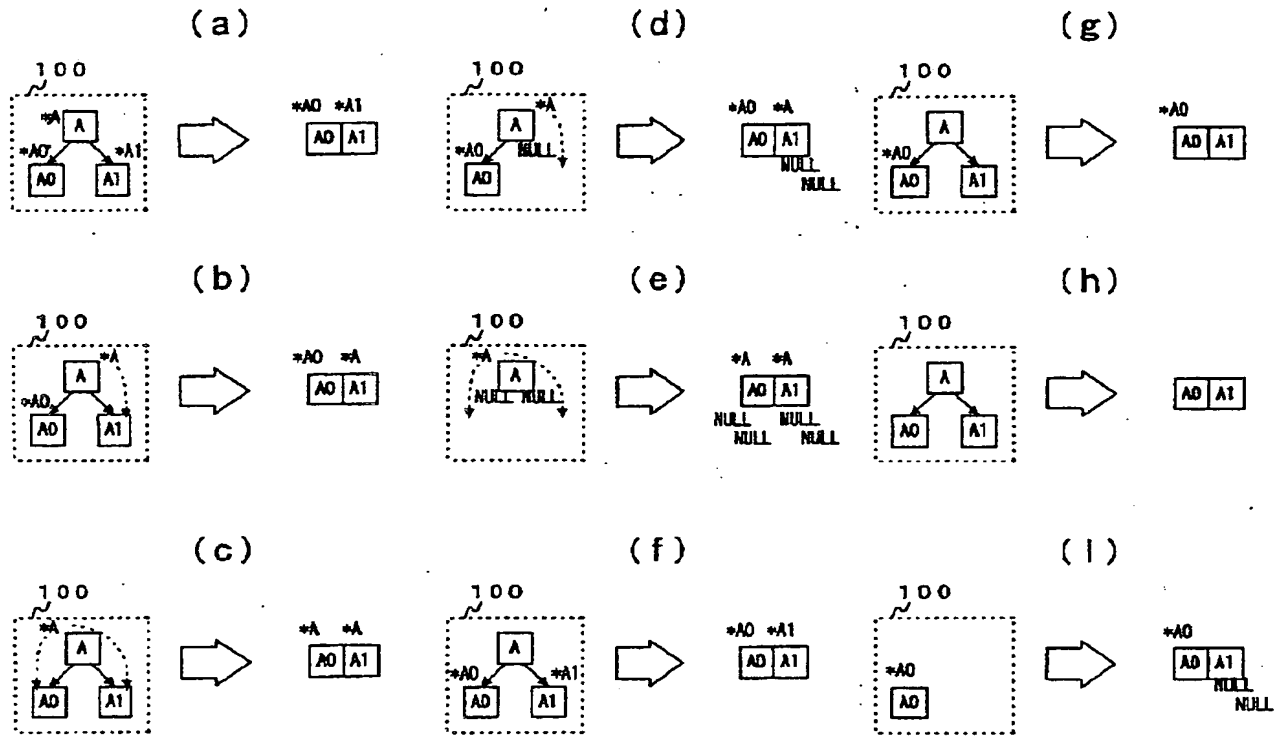


(b)

検査ビットの値	00	01	10	11
120	120			
121		121		
122			122	
123				123
124	124	124	124	124
125	125	125	125	125

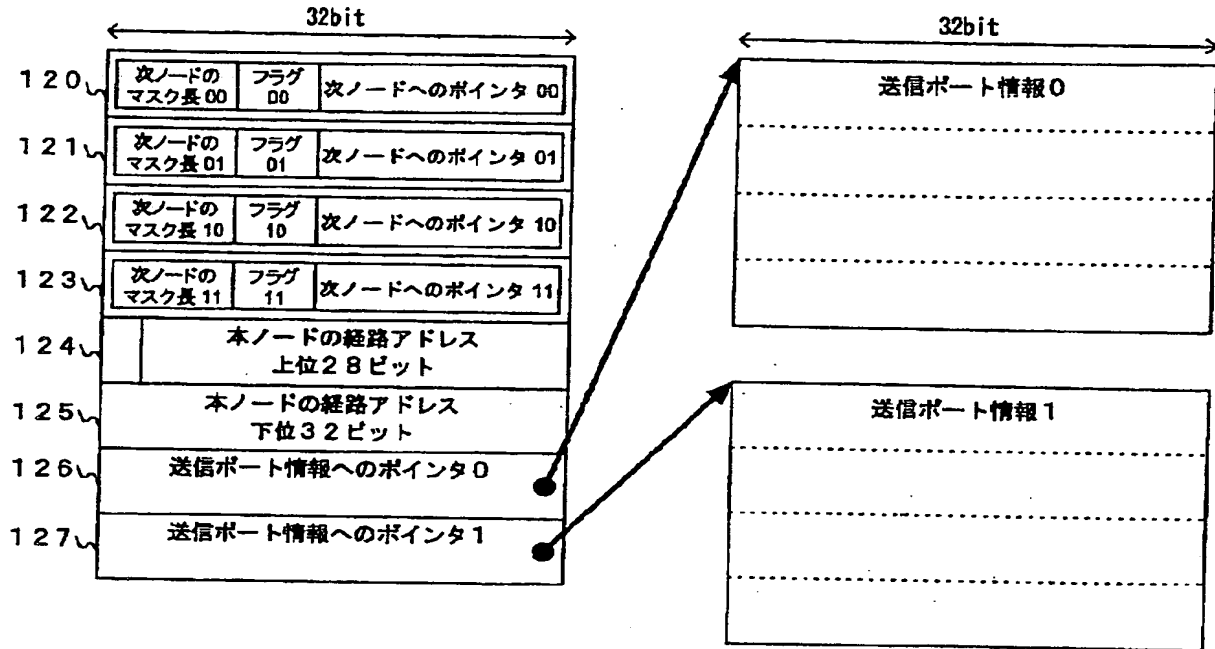
【図17】

図17



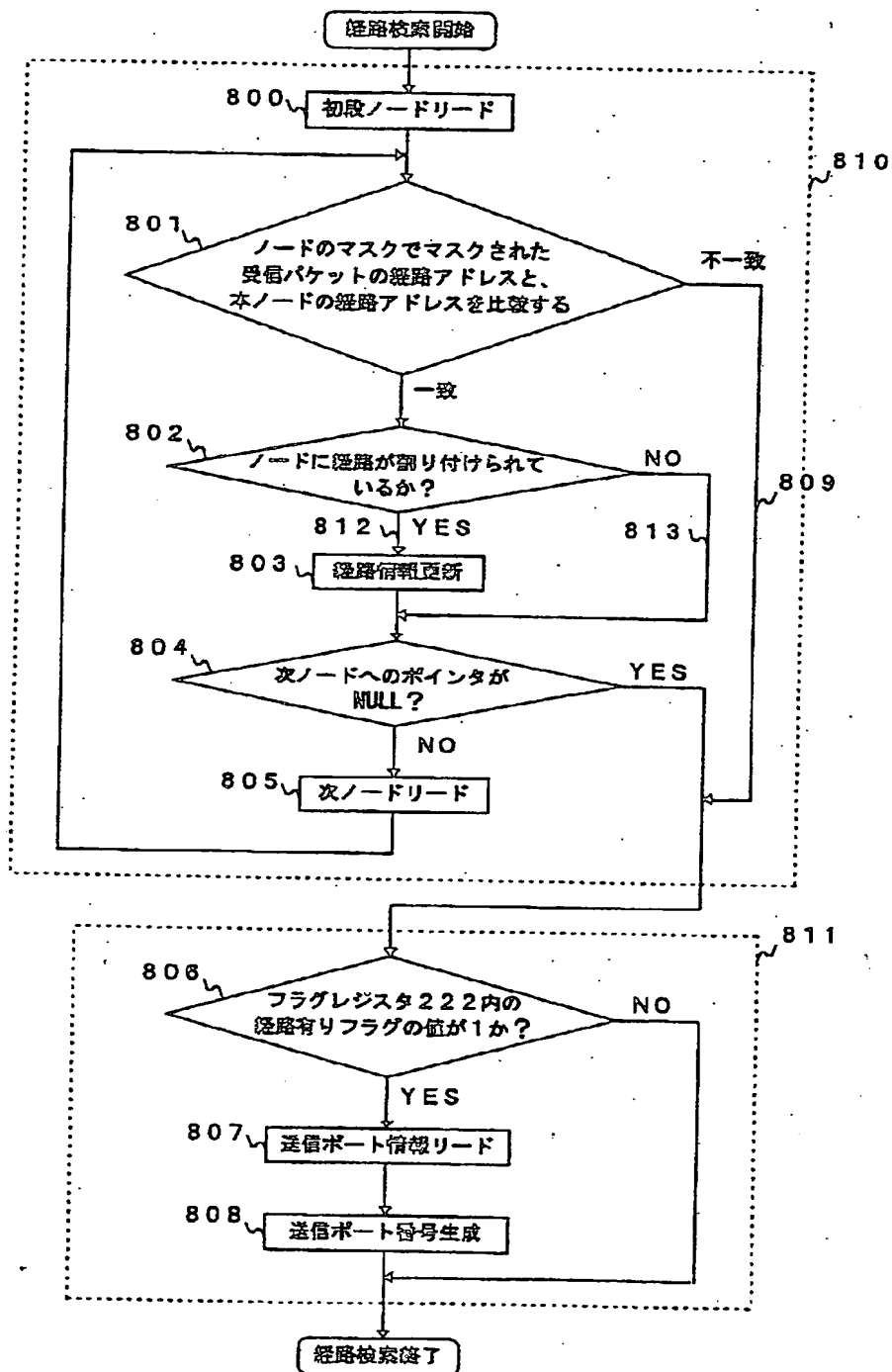
【図19】

図19



【図22】

図22



【図23】

図23

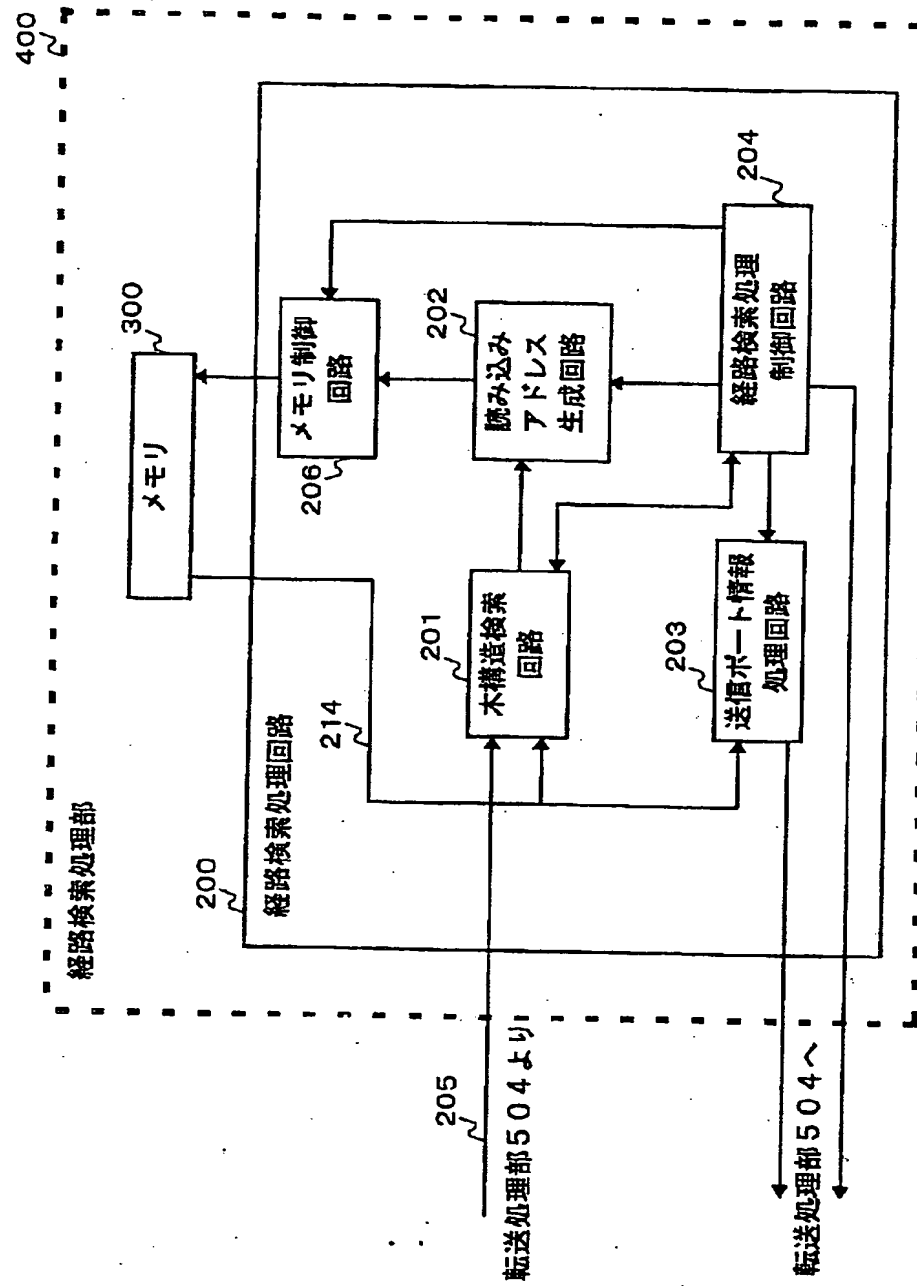
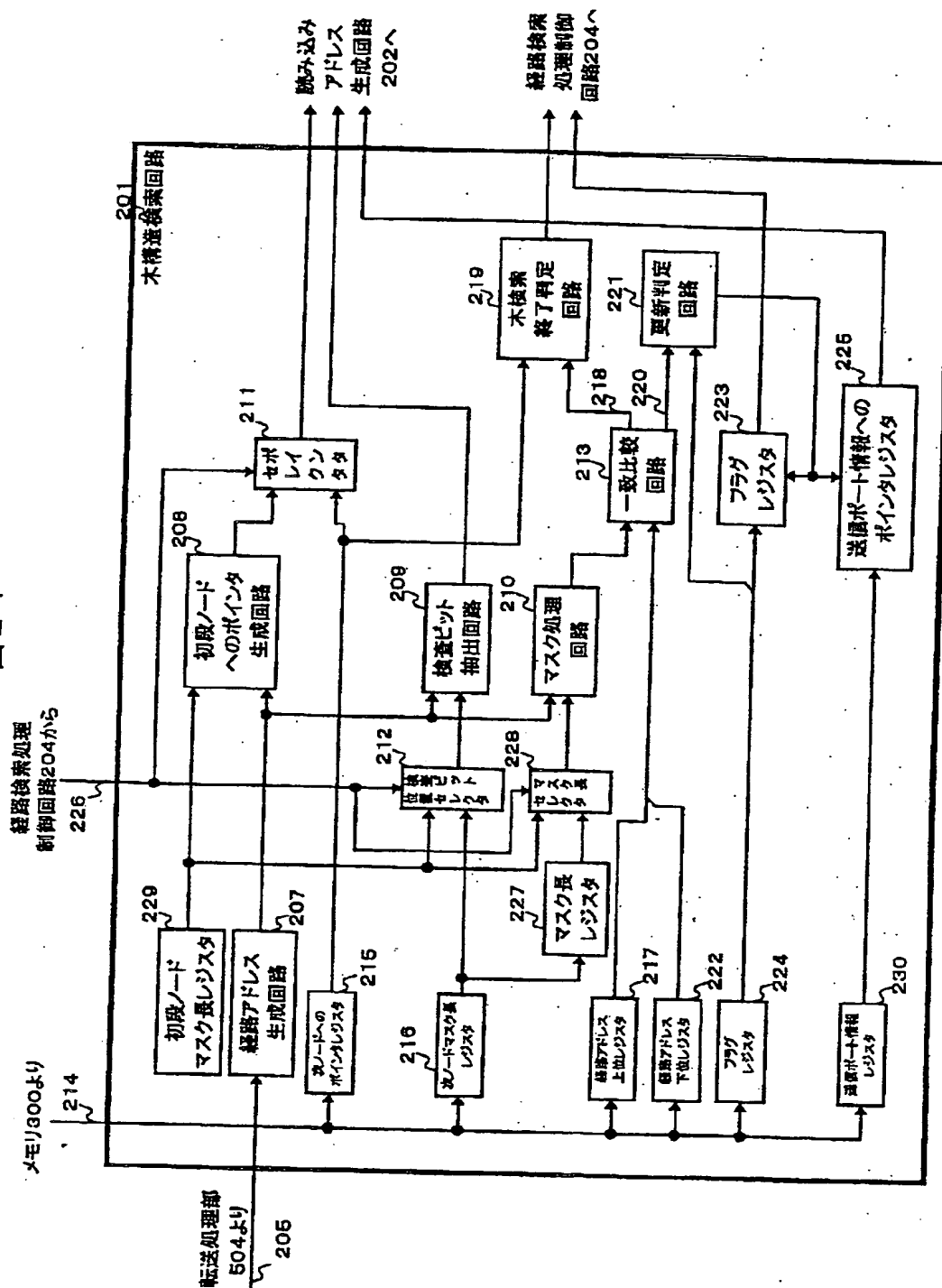


图 24



テーマコード（参考）

(72)発明者 須貝 和雄
神奈川県秦野市堀山下1番地 株式会社日立製作所汎用コンピュータ事業部内

(72)発明者 相本 毅
東京都国分寺市東恋ヶ窪一丁目280番地
株式会社日立製作所中央研究所内

(72)発明者 松山 信仁
神奈川県秦野市堀山下1番地 株式会社日立インフォメーションテクノロジー内

(72)発明者 左古 義人
神奈川県秦野市堀山下1番地 株式会社日立製作所汎用コンピュータ事業部内

(72)発明者 関野 浩
神奈川県秦野市堀山下1番地 株式会社日立インフォメーションテクノロジー内

Fターム(参考) 5B089 GA31 HB02 KA05 KB06 KC15
KC20 KC44 KE07 KG08
5K030 GA01 HA08 HB11 HB28 HC01
HC14 HD03 HD06 HD09 KA05
KX28 LB06 LD01
5K033 AA02 CB01 CB08 CB11 CB13
CC02 DA01 DA05 DB18 EC03
5K034 AA02 BB07 DD03 EE09 FF08
FF11 FF13 KK21 LL01